

Multidimensional data clustering and dimension reduction for indexing and searching

Patent Number: ☐ [US6122628](#)
Publication date: 2000-09-19
Inventor(s): CASTELLI VITTORIO (US); THOMASIAN ALEXANDER (US); LI CHUNG-SHENG (US)
Applicant(s): IBM (US)
Requested Patent: ☐ [JP11242674](#)
Application Number: US19970960540 19971031
Priority Number (s): US19970960540 19971031
IPC Classification: G06F17/30
EC Classification: [G06F17/30A](#)
Equivalents: CN1216841, DE69802960D, DE69802960T, ☐ [EP1025514](#) (WO9923578), [B1](#), HU0100581, JP3113861B2, PL340039, TW410304, ☐ [WO9923578](#)

Abstract

An improved multidimensional data indexing technique that generates compact indexes such that most or all of the index can reside in main memory at any time. During the clustering and dimensionality reduction, clustering information and dimensionality reduction information are generated for use in a subsequent search phase. The indexing technique can be effective even in the presence of variables which are not highly correlated. Other features provide for efficiently performing exact and nearest neighbor searches using the clustering information and dimensionality reduction information. One example of the dimensionality reduction uses a singular value decomposition technique. The method can also be recursively applied to each of the reduced-dimensionality clusters. The dimensionality reduction can also be applied to the entire database as a first step of the index generation.

Data supplied from the esp@cenet database - I2

【特許請求の範囲】

【請求項 1】多次元データを表示する方法であって、

(a) 前記多次元データを 1 つ以上のクラスタに区分するステップと、(b) 前記 1 つ以上のクラスタに対するクラスタ化情報を生成及び記憶するステップと、(c) 前記 1 つ以上のクラスタに対する 1 つ以上の次元縮小済みクラスタ及び次元縮小情報を生成するステップと、(d) 前記次元縮小情報を記憶するステップとを含んでいる、前記方法。

【請求項 2】前記 1 つ以上の次元縮小済みクラスタに対する次元縮小済み索引を生成及び記憶するステップを更に含んでいる、請求項 1 記載の方法。

【請求項 3】前記データが、複数のデータ・レコードを含んでいる空間データベース又はマルチメディア・データベース内に記憶され、

索引付けを遂行すべき前記データベースの表示を複数ベクトルの集合として作成するステップを更に含み、前記各ベクトルが、前記データベース内の 1 行に対応し、前記各ベクトルの要素が、当該ベクトルに対応する特定の行の列内に保持されている値に対応し、これらの列について、探索可能な索引が生成され、前記ステップ (a) が、前記ベクトルを 1 つ以上のクラスタに区分することを含んでいる、請求項 1 記載の方法。

【請求項 4】前記索引を、全体として計算機の主メモリ内に記憶するステップを更に含んでいる、請求項 2 記載の方法。

【請求項 5】前記ステップ (c) が、特異値分解を含み、

前記各クラスタに対する変換行列及び当該変換行列の固有値を生成するステップと、最大の固有値を含んでいる前記固有値の部分集合を選択するステップとを更に含み、前記次元縮小情報が、前記変換行列及び前記固有値の部分集合を含んでいる、請求項 2 記載の方法。

【請求項 6】前記次元縮小済み索引を使用して、指定データに最も類似する k 個のレコードを探索するために、前記記憶済みクラスタ化情報に基づいて、前記指定データを前記 1 つ以上のクラスタに関連付けるステップと、前記関連付けられたクラスタに対する前記記憶済み次元縮小情報に基づいて、前記指定データを当該クラスタに対する部分空間に射影するステップと、前記射影するステップに回答して、前記射影済み指定データに対する直交補空間を含んでいる次元縮小情報を生成するステップと、前記索引を介して、前記射影済み指定データに最も類似する k 個のレコードを有する前記関連付けられたクラスタを探索するステップと、関連付けられた他の任意のクラスタが前記射影済み指定データに最も類似する k 個のレコードのうち任意のレコ

ードを含み得るか否かを決定するステップと、前記射影済み指定データに最も類似する k 個のレコードのうち任意のレコードを含み得る前記任意のクラスタ上で、前記関連付けられたクラスタを探索するステップを反復するステップとを含んでいる、請求項 5 記載の方法。

【請求項 7】前記索引を介して、前記関連付けられたクラスタを探索するステップが、前記関連付けられたクラスタ内の k 最近傍と前記射影済み指定データとの間の距離 D を、ミスマッチ索引 δ^2 の関数として、

δ^2 (テンプレート, 要素) = D^2 (射影済みテンプレート, 要素) + $\Sigma \| \text{直交補空間} \|^2$

の如く計算するステップを含んでいる、請求項 6 記載の方法。

【請求項 8】前記指定データが、探索テンプレートを含み、

前記射影するステップが、前記次元縮小情報を使用して、前記テンプレートをそれが属するクラスタに関連する部分空間に射影するステップを含み、

前記射影済みテンプレートに対するテンプレート次元縮小情報を生成するステップを更に含み、

前記索引を介して、前記関連付けられたクラスタを探索するステップが、前記射影済みテンプレート及び前記テンプレート次元縮小情報に基づいて遂行され、

前記探索テンプレートに最も類似する k 個のレコードの k 最近傍集合を更新するステップを更に含んでいる、請求項 6 記載の方法。

【請求項 9】前記固有値の部分集合を選択するステップが、戻される結果の精度及び再現度の関数である、請求項 5 記載の方法。

【請求項 10】指定データに最も類似する k 個のレコードを探索するために、

前記クラスタ化情報に基づいて、前記指定データが属するクラスタを識別するステップと、

前記識別されたクラスタに対する前記次元縮小情報に基づいて、前記指定データの次元を縮小するステップと、

前記縮小するステップに回答して、次元縮小済み指定データに対する次元縮小情報を生成するステップと、

前記次元縮小情報を使用して、前記指定データが属するクラスタの次元縮小バージョンに対する多次元索引を探索するステップと、

前記多次元索引を介して、前記クラスタ内で前記最も類似する k 個のレコードを検索するステップと、

前記検索済みの最も類似する k 個のレコードのうち最も近いレコードよりも前記指定データに近いレコードを保持し得る他のクラスタを識別するステップと、

前記識別するステップに回答して、前記指定データに最も近い他の候補クラスタを探索するステップと、

前記他の候補クラスタの全てについて前記他のクラスタを識別するステップ及び前記他の候補クラスタを探索す

るステップを反復するステップとを更に含んでいる、請求項 2 記載の方法。

【請求項 1 1】前記次元縮小バージョンのクラスタ内の k 最近傍と前記射影済み指定データとの間の距離 D を、ミスマッチ索引 δ^2 の関数として、 δ^2 (テンプレート, 要素) = D^2 (射影済みテンプレート, 要素) + Σ 直交補空間 \parallel^2 の如く計算するステップを含んでいる、請求項 1 0 記載の方法。

【請求項 1 2】前記クラスタ化情報が、前記 1 つ以上のクラスタの重心に関する情報を含み、前記重心を一意的なラベルと関連付けるステップを更に含んでいる、請求項 1 記載の方法。

【請求項 1 3】前記データの次元が 8 より大きい、請求項 1 記載の方法。

【請求項 1 4】絶対探索を遂行するために、前記記憶済みクラスタ化情報に基づいて、指定データを 1 つの前記クラスタに関連付けるステップと、前記関連付けるステップに回答して、前記クラスタの次元縮小バージョンに対する記憶済み次元縮小情報に基づいて、前記指定データの次元を縮小するステップと、前記次元縮小済み指定データに基づいて、前記指定データにマッチする前記クラスタの次元縮小バージョンを探索するステップとを更に含んでいる、請求項 1 記載の方法。

【請求項 1 5】前記探索するステップが、線形走査を遂行することを含んでいる、請求項 1 4 記載の方法。

【請求項 1 6】前記ステップ (a) 乃至 (d) を再帰的に適用することにより、前記次元縮小済みクラスタの階層を作成するステップと、前記階層の最下位レベルにおけるクラスタに対する 1 つ以上の低次元索引を生成及び記憶するステップとを更に含んでいる、請求項 1 記載の方法。

【請求項 1 7】絶対探索を遂行するために、(1) 前記記憶済みクラスタ化情報を使用して、指定データが属するクラスタを探索すること、(2) 前記記憶済み次元縮小情報を使用して、前記階層の対応する最下位レベルに到達するまで、前記指定データの次元を縮小すること、及び (3) 前記低次元索引を使用して、前記指定データにマッチする前記クラスタの次元縮小バージョンを探索することを再帰的に適用するステップを更に含んでいる、請求項 1 6 記載の方法。

【請求項 1 8】類似性に基づく探索を遂行するために、(1) 前記記憶済みクラスタ化情報を使用して、指定データが属する前記クラスタを探索すること、及び (2) 前記記憶済み次元縮小情報を使用して、前記階層の最下位レベルに対応するように、前記指定データの次元を縮小することを再帰的に適用するステップと、前記指定データが属する前記階層の最下位レベルにある終端クラスタから開始して、前記階層の各レベルにおい

て、前記指定データの k 最近傍のうち 1 つ以上のものを保持し得る候補終端クラスタを探索するステップと、前記候補終端クラスタの各々ごとに、前記指定データに対する k 最近傍についてクラスタ内探索を遂行するステップとを更に含んでいる、請求項 1 6 記載の方法。

【請求項 1 9】類似性に基づく探索を遂行するために、前記指定データの次元を縮小するステップと、(1) 前記記憶済みクラスタ化情報を使用して、次元縮小済み指定データが属する前記クラスタを探索すること、及び (2) 前記記憶済み次元縮小情報を使用して、前記階層の最下位レベルに対応するように、前記次元縮小済み指定データの次元を縮小することを再帰的に適用するステップと、前記指定データが属する前記階層の最下位レベルにある終端クラスタから開始して、前記階層の各レベルにおいて、前記次元縮小済み指定データの k 最近傍のうち 1 つ以上のものを保持し得る候補終端クラスタを探索するステップと、前記候補終端クラスタの各々ごとに、前記次元縮小済み指定データに対する k 最近傍についてクラスタ内探索を遂行するステップとを更に含んでいる、請求項 1 6 記載の方法。

【請求項 2 0】前記データが、データベース内に記憶され、前記データベースの次元を縮小し且つ前記データベースに関連する次元縮小情報を生成するステップと、前記データベースに関連する前記次元縮小情報を記憶するステップとを更に含み、前記ステップ (a) が、前記次元縮小情報を生成するステップに回答して遂行される、請求項 1 記載の方法。

【請求項 2 1】絶対探索を遂行するために、前記データベースに対する前記次元縮小情報に基づいて、指定データの次元を縮小するステップと、前記縮小するステップに回答して、前記クラスタ化情報に基づいて、前記次元縮小済み指定データを 1 つの前記クラスタに関連付けるステップと、前記関連付けられたクラスタに対する前記次元縮小情報に基づいて、前記指定データの次元を、関連付けられたクラスタによって定義される前記次元縮小済みクラスタの次元に縮小するステップと、前記指定データの次元縮小済みバージョンに基づいて、これにマッチする次元縮小済みクラスタを探索するステップとを更に含んでいる、請求項 2 0 記載の方法。

【請求項 2 2】類似性に基づく探索を遂行するために、前記データベースに関連する前記次元縮小情報を使用して、指定データの次元を縮小するステップと、前記クラスタ化情報を使用して、前記次元縮小済み指定データが属する前記クラスタを識別するステップと、前記識別済みクラスタに対する前記次元縮小情報に基づいて、前記次元縮小済み指定データの次元を更に縮小す

るステップと、

次元を更に縮小した前記次元縮小済み指定データが属する前記クラスタの次元縮小バージョンを探索するステップと、

前記多次元索引を介して、前記クラスタ内の次元を更に縮小した前記次元縮小済み指定データに最も類似するk個のレコードを検索するステップと、

前記検索済みの最も類似するk個のレコードのうち最も遠いレコードよりも前記指定データに近いレコードを他のクラスタが保持し得るか否かを評価するステップと、
前記評価するステップに応答して、前記指定データに最も近い他のクラスタを探索するステップと、

前記他のクラスタの全てについて前記評価するステップ及び前記他のクラスタを探索するステップを反復するステップとを更に含んでいる、請求項20記載の方法。

【請求項23】前記データが、データベース内に記憶され、

前記1つ以上の次元縮小済みクラスタに対する1つ以上の次元縮小済み探索可能索引を生成及び記憶するステップを更に含んでいる、請求項20記載の方法。

【請求項24】絶対探索を遂行するために、前記記憶済みクラスタ化情報に基づいて、指定データを1つの前記クラスタに関連付けるステップと、
前記関連付けるステップに応答して、前記指定データを、前記関連付けられたクラスタ及び前記関連付けられたクラスタに対する前記記憶済み次元縮小情報によって定義される次元縮小済みクラスタに分解するステップと、

前記分解済み指定データに基づいて、これにマッチする前記次元縮小済みクラスタに対する前記索引を探索するステップとを更に含んでいる、請求項20記載の方法。

【請求項25】前記指定データが、探索テンプレートを含み、

前記関連付けるステップが、前記記憶済みクラスタ化情報に基づいて、前記探索テンプレートが属するクラスタを識別することを含み、

前記分解するステップが、前記記憶済みクラスタ化情報に基づいて、前記探索テンプレートを前記識別済みクラスタに対する部分空間に射影することを含み、

前記探索するステップが、前記射影済みテンプレートについてクラスタ内探索を遂行することを含んでいる、請求項24記載の方法。

【請求項26】(e)前記クラスタの形状に対する0次近似に対応するクラスタ境界を生成するステップと、

(f)極小外接ボックスによって前記各クラスタの形状を近似し且つそれから前記各クラスタの形状に対する1次近似を生成するステップと、(g)各次元の midpoint において前記外接ボックスを2k個の超方形に区分するステップと、(h)データ点を保持する超方形のみを保存し且つそれから前記クラスタの形状に対する2次近似を生

成するステップと、(i)前記保存済み超方形の各々ごとに、前記ステップ(g)及び(h)を反復することにより、前記クラスタに対する3次、4次、・・・、n次近似を逐次に生成するステップとを更に含んでいる、請求項1記載の方法。

【請求項27】各クラスタの形状構造に対する逐次近似の階層を探索するために、

前記データベースに関連する前記次元縮小情報を使用して、前記指定データの次元を縮小するステップと、

前記クラスタ化情報に基づいて、前記次元縮小済み指定データが属するクラスタを識別するステップと、

前記識別済みクラスタに対する前記次元縮小情報に基づいて、前記次元縮小済み指定データの次元を更に縮小するステップと、

次元を更に縮小した前記次元縮小済み指定データが属するクラスタの次元縮小済みバージョンを探索するステップと、

前記多次元索引を介して、前記クラスタ内の次元を更に縮小した前記次元縮小済み指定データに最も類似するk個のレコードを検索するステップと、

前記検索済みの最も類似するk個のレコードのうち最も遠いレコードよりも前記指定データに近いレコードを1つ以上の他のクラスタが保持し得るか否かを評価するステップと、

前記クラスタの境界に基づいて、前記他のクラスタが前記指定データのk個のレコードのうち任意のレコードを保持し得る場合にのみ、当該他のクラスタを保存するステップと、

前記形状に対し次第に細密となる近似に基づいて、前記保存済みクラスタが前記k最近傍のうち任意のものを保持し得るか否かを反復的に決定するとともに、当該クラスタが前記階層の最も細密なレベルにおいて受け入れられる場合にのみ前記保存済みクラスタを保存するステップと、

前記保存済みクラスタを保存するステップに応答して、前記保存済みクラスタを、前記データのk最近傍のうち1つ以上を保持する候補クラスタとして識別するステップとを更に含んでいる、請求項26記載の方法。

【請求項28】多次元データに対する1つ以上の次元縮小済み索引を含み、多次元データを表示する方法ステップを遂行するように計算機によって実行可能なプログラム命令を有形的に記憶しているプログラム記憶装置であって、

前記方法ステップが、(a)前記多次元データを1つ以上のクラスタに区分するステップと、(b)前記1つ以上のクラスタに対するクラスタ化情報を生成及び記憶するステップと、(c)前記1つ以上のクラスタに対する1つ以上の次元縮小済みクラスタ及び次元縮小情報を生成するステップと、(d)前記次元縮小情報を記憶するステップとを含んでいる、前記プログラム記憶装置。

【請求項 2 9】前記 1 つ以上の次元縮小済みクラスタに対する次元縮小済み索引を生成及び記憶するステップを更に含んでいる、請求項 2 8 記載のプログラム記憶装置。

【請求項 3 0】前記データが、複数のデータ・レコードを含んでいる空間データベース又はマルチメディア・データベース内に記憶され、索引付けを遂行すべき前記データベースの表示を複数ベクトルの集合として作成するステップを更に含み、前記各ベクトルが、前記データベース内の 1 行に対応し、前記各ベクトルの要素が、当該ベクトルに対応する特定の行の列内に保持されている値に対応し、これらの列について、探索可能な索引が生成され、前記ステップ (a) が、前記ベクトルを 1 つ以上のクラスタに区分することを含んでいる、請求項 2 8 記載のプログラム記憶装置。

【請求項 3 1】前記索引を、全体として計算機の主メモリ内に記憶するステップを更に含んでいる、請求項 2 9 記載のプログラム記憶装置。

【請求項 3 2】前記ステップ (c) が、特異値分解を含み、前記各クラスタに対する変換行列及び当該変換行列の固有値を生成するステップと、最大の固有値を含んでいる前記固有値の部分集合を選択するステップとを更に含み、前記次元縮小情報が、前記変換行列及び前記固有値の部分集合を含んでいる、請求項 2 9 記載のプログラム記憶装置。

【請求項 3 3】前記次元縮小済み索引を使用して、指定データに最も類似する k 個のレコードを探索するために、前記記憶済みクラスタ化情報に基づいて、前記指定データを前記 1 つ以上のクラスタに関連付けるステップと、前記関連付けられたクラスタに対する前記記憶済み次元縮小情報に基づいて、前記指定データを当該クラスタに対する部分空間に射影するステップと、前記射影するステップに回答して、前記射影済み指定データに対する直交補空間を含んでいる次元縮小情報を生成するステップと、前記索引を介して、前記射影済み指定データに最も類似する k 個のレコードを有する前記関連付けられたクラスタを探索するステップと、関連付けられた他の任意のクラスタが前記射影済み指定データに最も類似する k 個のレコードのうち任意のレコードを含み得るか否かを決定するステップと、前記射影済み指定データに最も類似する k 個のレコードのうち任意のレコードを含み得る前記任意のクラスタ上で、前記関連付けられたクラスタを探索するステップを反復するステップとを含んでいる、請求項 3 2 記載のプログラム記憶装置。

【請求項 3 4】前記索引を介して、前記関連付けられたクラスタを探索するステップが、前記関連付けられたクラスタ内の k 最近傍と前記射影済み指定データとの間の距離 D を、ミスマッチ索引 δ^2 の関数として、

$$\delta^2 (\text{テンプレート, 要素}) = D^2 (\text{射影済みテンプレート, 要素}) + \Sigma \parallel \text{直交補空間} \parallel^2$$

の如く計算するステップを含んでいる、請求項 3 3 記載のプログラム記憶装置。

10 【請求項 3 5】前記指定データが、探索テンプレートを含み、前記射影するステップが、前記次元縮小情報を使用して、前記テンプレートをそれが属するクラスタに関連する部分空間に射影するステップを含み、前記射影済みテンプレートに対するテンプレート次元縮小情報を生成するステップを更に含み、前記索引を介して、前記関連付けられたクラスタを探索するステップが、前記射影済みテンプレート及び前記テンプレート次元縮小情報に基づいて遂行され、

20 前記探索テンプレートに最も類似する k 個のレコードの k 最近傍集合を更新するステップを更に含んでいる、請求項 3 3 記載のプログラム記憶装置。

【請求項 3 6】前記固有値の部分集合を選択するステップが、戻される結果の精度及び再現度の関数である、請求項 3 2 記載のプログラム記憶装置。

【請求項 3 7】指定データに最も類似する k 個のレコードを探索するために、前記クラスタ化情報に基づいて、前記指定データが属するクラスタを識別するステップと、

30 前記識別されたクラスタに対する前記次元縮小情報に基づいて、前記指定データの次元を縮小するステップと、前記縮小するステップに回答して、次元縮小済み指定データに対する次元縮小情報を生成するステップと、前記次元縮小情報を使用して、前記指定データが属するクラスタの次元縮小バージョンに対する多次元索引を探索するステップと、

40 前記多次元索引を介して、前記クラスタ内で前記最も類似する k 個のレコードを検索するステップと、前記検索済みの最も類似する k 個のレコードのうち最も遠いレコードよりも前記指定データに近いレコードを保持し得る他のクラスタを識別するステップと、前記識別するステップに回答して、前記指定データに最も近い他の候補クラスタを探索するステップと、前記他の候補クラスタの全てについて前記他のクラスタを識別するステップ及び前記他の候補クラスタを探索するステップを反復するステップとを更に含んでいる、請求項 2 9 記載のプログラム記憶装置。

50 【請求項 3 8】前記次元縮小バージョンのクラスタ内の k 最近傍と前記射影済み指定データとの間の距離 D を、ミスマッチ索引 δ^2 の関数として、

δ^2 (テンプレート, 要素) = D^2 (射影済みテンプレート, 要素) + $\Sigma \parallel$ 直交補空間 \parallel^2

の如く計算するステップを含んでいる、請求項37記載のプログラム記憶装置。

【請求項39】前記クラスタ化情報が、前記1つ以上のクラスタの重心に関する情報を含み、

前記重心を一意的なラベルと関連付けるステップを更に含んでいる、請求項28記載のプログラム記憶装置。

【請求項40】前記データの次元が8より大きい、請求項28記載のプログラム記憶装置。

【請求項41】絶対探索を遂行するために、前記記憶済みクラスタ化情報に基づいて、指定データを1つの前記クラスタに関連付けるステップと、前記関連付けるステップに回答して、前記クラスタの次元縮小バージョンに対する記憶済み次元縮小情報に基づいて、前記指定データの次元を縮小するステップと、前記次元縮小済み指定データに基づいて、前記指定データにマッチする前記クラスタの次元縮小バージョンを探索するステップとを更に含んでいる、請求項28記載のプログラム記憶装置。

【請求項42】前記探索するステップが、線形走査を遂行することを含んでいる、請求項41記載のプログラム記憶装置。

【請求項43】前記ステップ(a)乃至(d)を再帰的に適用することにより、前記次元縮小済みクラスタの階層を作成するステップと、前記階層の最下位レベルにおけるクラスタに対する1つ以上の低次元索引を生成及び記憶するステップとを更に含んでいる、請求項28記載のプログラム記憶装置。

【請求項44】絶対探索を遂行するために、

(1) 前記記憶済みクラスタ化情報を使用して、指定データが属するクラスタを探索すること、(2) 前記記憶済み次元縮小情報を使用して、前記階層の対応する最下位レベルに到達するまで、前記指定データの次元を縮小すること、及び(3) 前記低次元索引を使用して、前記指定データにマッチする前記クラスタの次元縮小バージョンを探索することを再帰的に適用するステップを更に含んでいる、請求項43記載のプログラム記憶装置。

【請求項45】類似性に基づく探索を遂行するために、

(1) 前記記憶済みクラスタ化情報を使用して、指定データが属する前記クラスタを探索すること、及び(2) 前記記憶済み次元縮小情報を使用して、前記階層の最下位レベルに対応するように、前記指定データの次元を縮小することを再帰的に適用するステップと、前記指定データが属する前記階層の最下位レベルにある終端クラスタから開始して、前記階層の各レベルにおいて、前記指定データのk最近傍のうち1つ以上のものを保持し得る候補終端クラスタを探索するステップと、前記候補終端クラスタの各々ごとに、前記指定データに対するk最近傍についてクラスタ内探索を遂行するステ

ップとを更に含んでいる、請求項43記載のプログラム記憶装置。

【請求項46】類似性に基づく探索を遂行するために、前記指定データの次元を縮小するステップと、

(1) 前記記憶済みクラスタ化情報を使用して、次元縮小済み指定データが属する前記クラスタを探索すること、及び(2) 前記記憶済み次元縮小情報を使用して、前記階層の最下位レベルに対応するように、前記次元縮小済み指定データの次元を縮小することを再帰的に適用するステップと、

前記指定データが属する前記階層の最下位レベルにある終端クラスタから開始して、前記階層の各レベルにおいて、前記次元縮小済み指定データのk最近傍のうち1つ以上のものを保持し得る候補終端クラスタを探索するステップと、

前記候補終端クラスタの各々ごとに、前記次元縮小済み指定データに対するk最近傍についてクラスタ内探索を遂行するステップとを更に含んでいる、請求項43記載のプログラム記憶装置。

20 【請求項47】前記データが、データベース内に記憶され、

前記データベースの次元を縮小し且つ前記データベースに関連する次元縮小情報を生成するステップと、

前記データベースに関連する前記次元縮小情報を記憶するステップとを更に含み、

前記ステップ(a)が、前記次元縮小情報を生成するステップに回答して遂行される、請求項28記載のプログラム記憶装置。

【請求項48】絶対探索を遂行するために、

30 前記データベースに対する前記次元縮小情報に基づいて、指定データの次元を縮小するステップと、

前記縮小するステップに回答して、前記クラスタ化情報に基づいて、次元縮小済み指定データを1つの前記クラスタに関連付けるステップと、

前記関連付けられたクラスタに対する前記次元縮小情報に基づいて、前記指定データの次元を、関連付けられたクラスタによって定義される前記次元縮小済みクラスタの次元に縮小するステップと、

前記指定データの次元縮小済みバージョンに基づいて、これにマッチする次元縮小済みクラスタを探索するステップとを更に含んでいる、請求項47記載のプログラム記憶装置。

【請求項49】類似性に基づく探索を遂行するために、前記データベースに関連する前記次元縮小情報を使用して、指定データの次元を縮小するステップと、

前記クラスタ化情報を使用して、前記次元縮小済み指定データが属する前記クラスタを識別するステップと、

前記識別済みクラスタに対する前記次元縮小情報に基づいて、前記次元縮小済み指定データの次元を更に縮小するステップと、

次元を更に縮小した前記次元縮小済み指定データが属する前記クラスタの次元縮小バージョンを探索するステップと、
前記多次元索引を介して、前記クラスタ内の次元を更に縮小した前記次元縮小済み指定データに最も類似するk個のレコードを検索するステップと、
前記検索済みの最も類似するk個のレコードのうち最も遠いレコードよりも前記指定データに近いレコードを他のクラスタが保持し得るか否かを評価するステップと、
前記評価するステップにตอบสนองして、前記指定データに最も近い他のクラスタを探索するステップと、
前記他のクラスタの全てについて前記評価するステップ及び前記他のクラスタを探索するステップを反復するステップとを更に含んでいる、請求項4記載のプログラム記憶装置。

【請求項50】前記データが、データベース内に記憶され、
前記1つ以上の前記次元縮小済みクラスタに対する1つ以上の次元縮小済み探索可能索引を生成及び記憶するステップを更に含んでいる、請求項4記載のプログラム記憶装置。

【請求項51】絶対探索を遂行するために、
前記記憶済みクラスタ化情報に基づいて、指定データを1つの前記クラスタに関連付けるステップと、
前記関連付けるステップにตอบสนองして、前記指定データを、前記関連付けられたクラスタ及び当該関連付けられたクラスタに対する記憶済み次元縮小情報によって定義される前記次元縮小済みクラスタに分解するステップと、
前記分解済み指定データに基づいて、これにマッチする前記次元縮小済みクラスタに対する前記索引を探索するステップとを更に含んでいる、請求項4記載のプログラム記憶装置。

【請求項52】前記指定データが、探索テンプレートを含み、
前記関連付けるステップが、前記記憶済みクラスタ化情報に基づいて、前記探索テンプレートが属するクラスタを識別することを含み、
前記分解するステップが、前記記憶済みクラスタ化情報に基づいて、前記探索テンプレートを前記識別済みクラスタに対する部分空間に射影することを含み、
前記探索するステップが、前記射影済みテンプレートについてクラスタ内探索を遂行することを含んでいる、請求項51記載のプログラム記憶装置。

【請求項53】(e)前記クラスタの形状に対する0次近似に対応するクラスタ境界を生成するステップと、

(f)極小外接ボックスによって各クラスタの形状を近似し且つそれから各クラスタの形状に対する1次近似を生成するステップと、(g)各次元の中心において前記外接ボックスを2k個の超方形に区分するステップと、

(h)データ点を保持する超方形のみを保存し且つそれから前記クラスタの形状に対する2次近似を生成するステップと、(i)前記保存済み超方形の各々ごとに、前記ステップ(g)及び(h)を反復することにより、前記クラスタに対する3次、4次、・・・、n次近似を逐次に生成するステップとを更に含んでいる、請求項28記載のプログラム記憶装置。

【請求項54】各クラスタの形状構造に対する逐次近似の階層を探索するために、

前記データベースに関連する前記次元縮小情報を使用して、前記指定データの次元を縮小するステップと、
前記クラスタ化情報に基づいて、前記次元縮小済み指定データが属するクラスタを識別するステップと、
前記識別済みクラスタに対する前記次元縮小情報に基づいて、前記次元縮小済み指定データの次元を更に縮小するステップと、

次元を更に縮小した前記次元縮小済み指定データが属するクラスタの次元縮小済みバージョンを探索するステップと、

前記多次元索引を介して、前記クラスタ内の次元を更に縮小した前記次元縮小済み指定データに最も類似するk個のレコードを検索するステップと、
前記検索済みの最も類似するk個のレコードのうち最も遠いレコードよりも前記指定データに近いレコードを1つ以上の他のクラスタが保持し得るか否かを評価するステップと、

前記クラスタの境界に基づいて、前記他のクラスタが前記指定データのk個のレコードのうち任意のレコードを保持し得る場合にのみ、当該他のクラスタを保存するステップと、

前記形状に対し次第に細密となる近似に基づいて、前記保存済みクラスタが前記k最近傍のうち任意のものを保持し得るか否かを反復的に決定するとともに、当該クラスタが前記階層の最も細密なレベルにおいて受け入れられる場合にのみ前記保存済みクラスタを保存するステップと、

前記保存済みクラスタを保存するステップにตอบสนองして、前記保存済みクラスタを、前記データのk最近傍のうち1つ以上を保持する候補クラスタとして識別するステップとを更に含んでいる、請求項53記載のプログラム記憶装置。

【請求項55】多次元データを表示するための計算機可読プログラム・コード手段を有形的に記憶する記憶媒体を備えた計算機用プログラム製品であって、
前記計算機可読プログラム・コード手段が、
前記多次元データを1つ以上のクラスタに区分するように計算機を制御するためのクラスタ化手段と、
前記クラスタ化手段に結合され、前記1つ以上のクラスタに対するクラスタ化情報を生成及び記憶するように計算機を制御するための手段と、

前記クラスタ化手段に結合され、前記 1 つ以上のクラスタに対する 1 つ以上の次元縮小済みクラスタ及び次元縮小情報を生成するように計算機を制御するための次元縮小手段と、

前記次元縮小手段に結合され、前記次元縮小情報を記憶するように計算機を制御するための手段と、

前記 1 つ以上の次元縮小済みクラスタに対する次元縮小済み索引を生成及び記憶するように計算機を制御するための手段を含んでいる、前記計算機用プログラム製品。

【発明の詳細な説明】

【0001】

【関連する出願】本発明の関連出願は、“Searching Multidimensional Indexes Using Associated Clustering and Dimension Reduction Information”と題する、1997年10月31日に出願した係属中の米国特許出願（本出願人の整理番号 Y O 9 9 7 3 6 8）である。本明細書では、この関連出願の内容を援用する。

【0002】

【発明の属する技術分野】本発明は、多次元データのコンパクト表示を生成及び探索することに係り、更に詳細に説明すれば、関連付けられたクラスタ化情報及び次元縮小情報を使用して、データベース・システム内にある多次元データのコンパクト索引表示を生成及び探索することに係る。

【0003】

【従来の技術】多次元索引付けは、空間データベースにとって基本的な技法であり、地図情報システム（GIS）や、膨大なデータ・ウェアハウスを使用して意思決定を支援するためのオンライン分析処理（OLAP）や、イメージ及びビデオ・データから高次元の特徴ベクトルを抽出するようなマルチメディア・データベースに対し広範に適用することができる。

【0004】意思決定支援システムは、事業体が成功するための不可欠の技術になりつつある。各事業体は、意思決定支援システムを使用して、商用データベースから（データ・ウェアハウスとも呼ばれる）有用な情報を推論することができる。商用データベースが状態情報を保持するのに対し、データ・ウェアハウスは履歴情報を保持するのが普通である。一般に、データ・ウェアハウスのユーザは、個々のレコードを別々に調べるのではなく、傾向を把握することに関心を持っている。従って、意思決定用の照会を主体とする計算を集中的に遂行して、集約機能を多量に使用する。その結果、計算が完了するまでの遅れが長くなって、生産性上の制約が受け入れがたいものになることがある。

【0005】かかる遅れを減少させるのに使用されていた公知の技法は、出現頻度が高い照会を予め計算するか、又はサンプリング技法を使用する、というものである。特に、最近の注目を集めているのは、データ・キューブのようなオンライン分析処理（OLAP）技法を、

大規模な関係データベース又は意思決定支援用のデータ・ウェアハウスに適用することである（例えば、Jim Gray et al, “Data Cube: A Relational Aggregation Operator Generalizing Group-By, Cross-Tab, and Sub-Totals”, International Conference on Data Engineering, 1996, New Orleans, pp. 152-160 を参照）。ここでは、ユーザは、データ・ウェアハウスからの履歴データを多次元データ・キューブとして観察するのが普通である。このキューブ内の各セル（又は格子点）は、総売上高のような諸項目の集約から成るビューである。

【0006】多次元索引付けは、異なるタイプの照会に回答するために使用することができる。これらの照会を例示すると、次の通りである。

(1) 索引付けした列の指定値を有するレコードを探索する（絶対探索）。

(2) $[a_1, \dots, a_n]$ 、 $[b_1, \dots, b_n]$ 、 \dots 、 $[z_1, \dots, z_n]$ 内のレコードを探索する（範囲探索）。但し、 a 、 b 及び z は、異なる次元を表すものとする。

(3) ユーザが指定したテンプレート又は例に対し最も類似する k 個のレコードを探索する（ k 最近傍探索）。

【0007】多次元索引付けは、イメージ・マイニングにも適用することができる。イメージ・マイニング製品の一例は、「MEDIAMINER」（本出願人の商標）である。この製品が提供する 2 つのツール、すなわち「イメージ・コンテンツ照会プログラム（QBIC）」及び「IMAGEMINER」は、手動的に作成した関連キーワードのリストを使用して探索を遂行するのではなく、コンテンツを分析することによってイメージの検索を遂行する。

【0008】QBIC を使用するのに適したアプリケーションは、キーワードを使用しても適切な結果が得られないようなもの、例えば博物館及び美術館用のライブラリに關係するアプリケーション、又は電子商取引のオンライン在庫品の写真に關係するアプリケーションである。ちなみに、後者のアプリケーションでは、顧客がビジュアル・カタログを使用し、必要とする商品（例えば、壁紙や衣服）のカラーやテクスチャなどを吟味した上で、所望の商品を探索することができる。

【0009】「IMAGEMINER」のようなイメージ・マイニング用アプリケーションは、概念的な照会（例えば、「森の景色」、「氷」、「シリンダ」など）を使用して、イメージ・データベースに対する照会を遂行することができる。カラーやテクスチャ及び輪郭などのイメージ・コンテンツは、システムが自動的に認識可能な単純オブジェクトとして結合される。

【0010】これらの単純オブジェクトは、知識ベース内に表示される。これを分析して得られるテキストの記述は、後の検索に備えて索引付けされる。

【0011】データベースの照会を実行中、データベー

10

20

30

40

50

ス探索プログラムは、記憶データの一部及び索引付け構造の一部をアクセスする。この場合、アクセスされるデータ量は、照会のタイプ、ユーザが提供するデータ及び索引付けアルゴリズムの効率に依存する。大規模データベースでは、そのデータ及び索引付け構造の少なくとも一部を、計算機システムのメモリ階層のうち 1 台以上のハードディスクから成る大容量の低速部分に記憶するのが普通である。探索プロセスの間、かかるデータ及び索引付け構造の一部を、メモリ階層の高速部分、すなわち主メモリ及びキャッシュ・メモリの 1 つ以上のレベルにロードする。一般に、メモリ階層の高速部分は、比較的高価であるという理由で、メモリ階層の記憶容量のうち数パーセントに相当するにすぎない。キャッシュ・メモリの 1 つ以上のレベルに完全にロード可能な命令及びデータを使用するプログラムは、主メモリ内に存在する命令及びデータを使用するプロセスよりも高速及び効率的であり、また後者のプロセスは、ハードディスク上に存在する命令及びデータを使用するプログラムよりも高速である。この場合の技術的制限とは、主メモリ及びキャッシュ・メモリのコストが比較的高いために、大規模データベースを完全に記憶するのに十分な主メモリ及びキャッシュ・メモリを備えた計算機システムを構築することが実際のでない、というものである。

【0012】従って、当分野において要請されている索引付け技法とは、任意の時点で殆ど又は全ての索引を主メモリ内に常駐可能にするようなサイズの索引を生成するとともに、探索プロセス中にハードディスクから主メモリに転送すべきデータの量を制限するようなものである。本発明は、かかる要請を実現することに向けられている。

【0013】Rツリーを含む周知の空間索引付け技法は、これを範囲照会及び最近傍照会のために使用することができる。Rツリーの詳細は、例えば A. Gutman, "R-trees: A Dynamic Index Structure for Spatial Searching", ACM SIGMOD Conf. on Management of Data, Boston, MA, June 1994 に記述されている。しかしながら、これらの技法は、特徴空間の次元の数が増大するにつれて、探索空間が次第に疎になるために、その効率が急速に低下する。例えば、次元の数が 8 よりも大きい場合、Rツリーのような方法は有用でないことが知られている。この場合の有用基準は、要求を完了するための時間と、データベース内の全てのレコードを順次に走査してこの要求を完了するという暴力的戦略が必要とする時間との比較によって決まる。高次元空間における通常の索引付け技法の非効率性は、「次元の呪い」と呼ばれる現象に起因する。この現象については、例えば V. Cherkassky et al, "From Statistics to Neural Networks", NATO ASI Series, Vol. 136, Springer-Verlag, 1994 に記述されている。また、この現象に起因して、高次元の特徴空間については、索引空間を複数の超立方体

にクラスタ化することが非効率的となる。

【0014】前述のように、高次元の特徴空間を索引付けするための既存の空間索引付け技法が非効率的であることから、周知の技法は、特徴空間の次元の数を縮小することに向けられている。例えば、次元の数を縮小するために、(特徴選択とも呼ばれる) 可変部分集合選択を使用するか、又は特異値分解の後に可変部分集合選択を使用することができる(例えば、C. T. Chen, "Linear System Theory and Design", Holt, Rinehart and Winston, Appendix E, 1984 を参照)。可変部分集合選択は、統計分野で盛んに研究されていて、多数の方法が既に提案されている(例えば、Shibata et al, "An Optimal Selection of Regression Variables", Biometrika, Vol. 68, No. 1, 1981, pp. 45-54 を参照)。これらの方法が索引生成システムにおいて有効であるのは、多くの変数(データベース内の列)が高度に相関している場合だけである。

【0015】従って、高度に相関していない変数の存在下であっても、高次元データの索引付けを効率的に遂行するような改良された索引付け技法が要請されている。この技法は、メモリの利用度及び探索速度の観点から、索引を効率的に生成しなければならない。本発明は、これらの要請を実現することに向けられている。

【0016】

【発明が解決しようとする課題】前述の要請に従って、本発明の目的は、多次元データのコンパクト表示を生成するための装置及び方法を提供することにある。本発明は、データベース用の探索可能な多次元索引を生成することを特徴としている。また、本発明は、これらの索引を柔軟に生成するとともに、絶対探索及び類似性に基づく探索を効率的に遂行することを特徴としている。更に、本発明は、探索プロセス中にハードディスクから主メモリに転送するデータの量を制限することを可能にする、コンパクト索引を生成することを特徴としている。

【0017】本発明の 1 つの適用例は、多次元索引付けである。多次元索引付けは、空間データベースにとって基本的な技法であり、地図情報システム(GIS)や、膨大なデータ・ウェアハウスを使用して意思決定を支援するためのオンライン分析処理(OLAP)や、イメージ及びビデオ・データから高次元の特徴ベクトルを抽出するマルチメディア・データベースのマイニングを遂行するためのイメージ・マイニング製品に対し広範に適用することができる。

【0018】

【課題を解決するための手段】本発明の特徴を有する方法の一例は、(a) 多次元データを 1 つ以上のクラスタに区分するステップと、(b) 前記クラスタに対するクラスタ化情報を生成及び記憶するステップと、(c) 前記クラスタに対する 1 つ以上の次元縮小済みクラスタ及び次元縮小情報を生成するステップと、(d) 前記次元

縮小情報を記憶するステップとを含んでいる。また、本発明は、前記次元縮小済みクラスタに対する次元縮小済み索引を生成及び記憶することを特徴としている。

【0019】各クラスタ内で使用されている索引付け技法に依存して、対応する索引付け機構を使用することにより、目標ベクトルを検索することができる。例えば、各クラスタ内の索引付けのために、通常の多次元空間索引付け技法（Rツリーを含むがこれに限らない）を使用することができる。代替的に、どんな空間索引付け構造も利用できない場合は、クラスタ内探索機構として、暴力的又は線形走査を利用することができる。

【0020】推奨実施例において、ステップ（c）は、特異値分解であり、かくて分解済みの指定データに基づいて、次元縮小済みクラスタにマッチする索引を探索する。次元縮小情報の一例は、（固有値及び固有ベクトルを含む）変換行列である。この変換行列は、特異値分解及びこの変換行列の選択された固有値によって生成される。

【0021】本発明の特徴に従った多次元索引を生成する方法の他の例は、索引付けを遂行すべきデータベースの表示を、複数ベクトルの集合として作成するステップを含んでいる。但し、各ベクトルは、このデータベース内の1行に対応し、各ベクトルの要素は、このベクトルに対応する特定の行の列（そのための索引を生成しなければならない）内に保持されている値に対応する。この方法では、クラスタ化技法を使用して、このベクトル集合を（クラスタとも呼ばれる）1つ以上のグループに区分するとともに、このクラスタ化に関連するクラスタ化情報を生成及び記憶する。次いで、次元縮小情報を使用して、これらのグループの各々に対し次元縮小技法を別個に適用することにより、クラスタ内にある諸要素の低次元表示を生成する。次に、次元縮小済みクラスタの各々ごとに、このクラスタの次元の数について効率的な索引を生成する技法を使用して、索引を生成する。

【0022】本発明の他の特徴に従って、この方法は、これが再帰的となるように、次元縮小済みクラスタの各々に対し別個に適用することができる。次元縮小ステップ及びクラスタ化ステップの両者を再帰的に適用するというこのプロセスが終了するのは、もはや次元を縮小させることができないような場合である。

【0023】本発明の他の特徴に従って、次元縮小ステップは、（データベースを区分するステップの前に）索引生成プロセスにおける最初のステップとして、データベースの全体に適用することができる。（クラスタ化とも呼ばれる）区分ステップ及び次元縮小ステップの間、探索段階で使用するのに備えて、クラスタ化情報及び次元縮小情報を生成する。

【0024】本発明の更に他の特徴に従って、次元縮小ステップを容易にするように、クラスタ化ステップを適当に選択することができる。例えば、これを遂行するに

は、ユークリッド距離のような空間的に不変の距離関数から導かれる損失を最小にするのではなく、データの局所共分散構造に従って空間を区分するというクラスタ化方法を使用することができる。

【0025】また、本発明は、指定データに最も類似する検索済みのk個の要素のうち最も遠い要素よりも指定データに近い要素を、他のクラスタが保持し得るか否かを評価することを特徴としている。公知のように、クラスタ化情報を使用すると、複数の区分の境界を再構成することが可能であり、またこれらの境界を使用すると、指定データに最も類似するk個の要素のうち1つの要素を、1つのクラスタが保持し得るか否か決定することができる。当業者には明らかなように、これらのクラスタ境界は、クラスタ自体の構造に対する簡単な近似である。すなわち、この境界の数学的形式からは、この境界上の所与の位置にクラスタの要素が存在するか否かを断定することはできない。一例として、データベースが2つの球型データ・クラスタを保持しており、そしてこれらのクラスタが互いに著しく離れているようなケースを検討する。このケースの妥当な境界は、これらのクラスタの重心を結合する線分に垂直で且つこれらの重心から等距離の超平面となろう。これらのクラスタは互いに著しく離れているから、この境界近傍には如何なるデータ点も存在しない。他のケースでは、この境界は、両クラスタの多数の要素に非常に接近することがあり得る。従って、本発明は、各クラスタの実際の形状構造に対する諸近似の階層を使用することにより、このクラスタが、指定データに最も類似するk個の要素のうち1つ以上の要素を保持し得るか否かを決定することを特徴としている。

【0026】推奨実施例において、本発明は、計算機可読プログラム記憶媒体上に有形的に記憶したソフトウェアとして実現することができる。このソフトウェアを構成するプログラム命令を計算機によって実行すると、多次元データを表示する次の方法、すなわち多次元データのコンパクト表示を生成するステップと、データベース用の探索可能な多次元索引を生成するステップと、これらの索引を使用して絶対探索及び類似性に基づく探索を効率的に遂行するステップとを含む方法が実施される。

【0027】

【発明の実施の形態】図1には、本発明の特徴を実現したネットワーク化クライアント／サーバ・システムが例示されている。図示のように、複数クライアント101及び複数サーバ106は、ネットワーク102によって相互接続されている。サーバ106は、データベース管理システム（DBMS）104及び直接アクセス記憶装置（DASD）105を含んでいる。一般に、クライアント101上で照会を生成し、これをネットワーク102を通してサーバ106に発信する。かかる照会は、ユーザが提供する例又は探索テンプレートのような指定デ

ータを含んでおり、DASD 1 0 5 内に記憶されているデータベースを検索又は更新を遂行するためにDBMS 1 0 4 と対話する。DBMS 1 0 4 の一例は、IBM 社から提供されている「DB 2」（本出願人の商標）である。

【0 0 2 8】本発明の 1 つの側面に従って、多次元の、例えば空間索引付けを必要とする照会（範囲照会及び最近傍照会を含む）は、多次元索引付けエンジン 1 0 7 を呼び出す。この多次元索引付けエンジン 1 0 7（図 8 ～ 図 1 1 を参照して後述）は、本発明の索引生成論理 1 1 0（図 6 及び図 7 を参照して後述）が生成する 1 つ以上のコンパクト多次元索引 1 0 8、クラスタ化情報 1 1 1 及び次元縮小情報 1 1 2 に基づいて、この照会が指定する制約を満足するようなベクトル又はレコードを検索する。本発明のコンパクト索引 1 0 8 の全部又は殆どは、サーバ 1 0 6 の主メモリ及びキャッシュ・メモリの少なくとも一方に記憶することが好ましい。空間データベースのようなデータベースは、1 つ以上のシステム上に存在することができる。また、多次元索引付けエンジン 1 0 7 及び索引生成論理 1 1 0 の少なくとも一方は、DBMS 1 0 4 の一部として統合化することができる。更に、（探索論理とも呼ばれる）多次元索引付けエンジン 1 0 7 及び索引生成論理 1 1 0 は、サーバ 1 0 6 上で実行可能な計算機プログラム製品上のソフトウェアとして有形的に実施することができる。

【0 0 2 9】1 つの適用例は、スーパーマーケットの記憶済み販売時点情報管理トランザクションであり、店の位置（緯度及び経度）の幾何座標を含むものである。ここでは、サーバ 1 0 6 は、記憶済みデータから知識又はパターンを探索するために意思決定支援タイプのアプリケーションをサポートしていることが好ましい。例えば、オンライン分析処理（OLAP）エンジン 1 0 3 を使用して、OLAP に関係する照会をインターセプトすることにより、これらの処理を容易に遂行することができる。本発明に従って、OLAP エンジン 1 0 3 は、可能であれば DBMS 1 0 4 と関連して、多次元索引エンジン 1 0 7 を使用することにより、OLAP に関係する照会用の索引 1 0 8 を探索する。なお、本発明の索引生成論理 1 1 0 は、データ・ウェアハウスの多次元データ・キューブ表示にも適用することができる。データ・ウェアハウスの多次元データ・キューブ表示を生成するための方法及びシステムは、“System and Method for Generating Multi-Representation of a Data Cube”と題する、1 9 9 7 年 4 月 1 4 日に出願された係属中の米国特許出願第 8 4 3 2 9 0 号に記述されている。本明細書では、この出願の内容を援用する。

【0 0 3 0】空間索引付けの効果を享受するデータの他の例は、マルチメディア・データである。オーディオ、ビデオ及びイメージのようなマルチメディア・データは、索引付けのために使用するメタデータとは別個に記

憶することができる。メディア・データの索引付け及び検索を容易にするために使用可能なメタデータの 1 つの重要な成分は、生データから抽出されるような「特徴ベクトル」である。例えば、イメージの領域からテクスチャ、カラー・ヒストグラム及び形状を抽出し、これを検索用の索引 1 0 8 を構成するために使用することができる。

【0 0 3 1】イメージ・マイニング・アプリケーションの例は、IBM 社の「DB2 Image Extender」における統合化探索機構の Q B I C である。Q B I C は、イメージ照会エンジン（サーバ）、並びに HTML グラフィカル・ユーザ・インタフェース及び関連する共通ゲートウェイ・インタフェース（C G I）スクリプトから成るサンプル・クライアントを含んでおり、これらが相まって完全なアプリケーションの基礎を形成している。このサーバ及びクライアントの両者は、ユーザがアプリケーションに特有のイメージ・マッチング関数を開発し且つこれを Q B I C に追加することができるように、拡張可能となっている。また、イメージ探索サーバは、ビジュアル・イメージのコンテンツに基づいて、大規模イメージ・データベースの照会を可能にする。

【0 0 3 2】この特徴は、次の通りである。

（1）ビジュアル・メディアの形式で照会すること。例えば、「これと同様のイメージを示せ」という具合。但し、「同様」の意味を、カラーや、レイアウト、テクスチャなどで定義する必要がある。

（2）照会イメージに対する類似性に従って、イメージをランク付けすること。

（3）諸イメージの自動的索引付けを遂行すること。但し、カラー及びテクスチャの数値記述子を記憶する必要がある。探索の間、類似イメージを探索するために、これらのプロパティを使用する。

（4）ビジュアル照会とテキスト照会又は日付のようなパラメータに関する照会との組み合わせ。

【0 0 3 3】同様に、索引付けを遂行すべきデータベースの表示を複数ベクトルの集合として作成することにより、複数の索引を生成することができる。但し、各ベクトルは、このデータベース内の 1 行に対応し、各ベクトルの要素は、この特定の行について、それぞれの列（そのための索引を生成することが必要である）に保持されている値に対応する。

【0 0 3 4】データベースの表示を複数ベクトルの集合として作成することは、当分野では周知である。この表示を作成するための代表的な方法は、データベースの各行ごとに、生成すべき索引の次元に等しい長さを有する 1 つの配列を作成し、この配列の諸要素に対し、この配列に対応するデータベースの特定の行の複数の列（そのための索引を生成しなければならない）に保持されている値をコピーする、というものである。

【0 0 3 5】ここで、ベクトル v の i 番目の要素が v_i

であると仮定すると、ベクトル v を次のように表すことができる。

【0036】

【数1】

$$v = [v_1 \dots v_N] \quad (1)$$

【0037】但し、 N は、索引付けのために使用するベクトルの次元の数である。

【0038】一般に、クライアント側は、3つのタイプの照会を指定することができる。これらの照会の全ては、所定の形式の空間索引付けを必要とする（従来技術の項の説明を参照）。これらの照会は、次の通りである。

(1) 絶対照会：1つのベクトルを指定すると、このベクトルにマッチするレコード又はマルチメディア・データが検索される。

(2) 範囲照会：このベクトルの各次元の下限及び上限が指定される。

(3) 最近傍照会：類似性の測度に基づいて、最も「類似する」ベクトルが検索される。

【0039】2つのベクトル v_1 及び v_2 の間の最も普通に使用されている類似性の測度は、ユークリッド距離の測度 d であり、これは次のように定義されている。

【0040】

【数2】

$$d^2 = \sum (v_{1,i} - v_{2,i})^2 \quad (2)$$

【0041】ここで、次元 i の全てが範囲照会又は最近傍照会の計算に必ずしも関与する必要はないことに留意されたい。両ケースにおいて、結果を検索するために、これらの次元の部分集合を指定することができる。

【0042】図2には、多次元空間における複数のベクトルの分布が例示されている。図示のように、この空間の全体を表示するには、3次元が必要である。しかしながら、 $x-y$ 平面、 $y-z$ 平面及び $z-x$ 平面上に位置するクラスタ 201、202 及び 203 の各々を表示するには、2次元だけが必要である。従って、かかるデータの適正なクラスタ化を通して、次元縮小を達成可能であると結論することができる。これと同じ次元縮小は、特異値分解だけでは達成することができない。なぜなら、特異値分解は、特徴空間内の軸が主要な次元（この例では3）とマッチするように、特徴空間を再配向するにすぎないからである。

【0043】1つのベクトルの1つ以上の次元を除去することは、元の点を部分空間に射影することと等価である。式(2)は、このベクトル内の個々の要素が異なっているような次元だけを計算すればよいことを示してい

る。その結果、このベクトルを部分空間に射影しても、距離の計算には影響しないから、除去される要素は元の空間内で変化しないことになる。

【0044】図3には、元の3次元空間内の3点を2次元部分空間に射影するに際し、これらの3点のうち任意の2点間の相対距離を維持するようにした距離計算の一例が示されている。図示のように、元の3次元空間では、点301と点302との間の距離は、点301と点303との間の距離よりも大きい。ここで、これらの点を2次元部分空間に射影した結果である点(304、305、306)間の相対距離が維持されていることに留意されたい。

【0045】図4には、3次元空間内の3点を2次元部分空間に射影するに際し、相対距離のランクに影響を及ぼすようにした一例が示されている。図示のように、3次元空間内の点401と点402との間の距離は、点402と点403との間の距離よりも大きい。しかしながら、このケースでは、射影済みの点404と点405との間の距離は、射影済みの点405と点406との間の距離よりも小さい。従って、射影済みの2次元部分空間では、2点間の相対距離が維持されないことになる。

【0046】以下では、複数のベクトルを部分空間に射影する際に生じ得るような最大誤差を評価するための技法を導出する。まず、このプロセスは、最大誤差の上界を決定することから開始する。1つのクラスタの重心 V_c を、次のように定義する。

【0047】

【数3】

$$V_c = \frac{1}{N} \sum V_i \quad (3)$$

【0048】但し、 N は、このクラスタ内にある複数のベクトル $\{V_1, \dots, V_N\}$ の総数である。このクラスタを k 次元の部分空間に射影した後、そこで一般性を失わないように最後の $(n-k)$ 次元を除去するものとする、元の空間に比較して、この部分空間における任意の2ベクトル間の距離に対する誤差が生ずることになる。この誤差の項は、次のように表される。

【0049】

【数4】

$$Error^2 = \sum_{i=k+1}^n (V_{1,i} - V_{2,i})^2 \quad (4)$$

【0050】式(4)から、次の不等式が直ちに成立する。

【0051】

【数5】

$$\begin{aligned}
 Error^2 &\leq \sum_{i=k+1}^n (|V_{1,i}| + |V_{2,i}|)^2 \\
 &\leq \sum_{i=k+1}^n (2 \max(|V_{1,i}|, |V_{2,i}|))^2 \\
 &\leq 4 \sum_{i=k+1}^n \max(|V_{1,i}|, |V_{2,i}|)^2 \quad (5)
 \end{aligned}$$

【0052】式(5)は、射影済みの部分空間における距離を計算する際に生ずる最大誤差が束縛されていることを示している。

【0053】図5には、本発明に従った距離の計算を遂行する際の近似の一例が示されている。テンプレート点 T (501) と生成点 V (506) との間の距離は、前掲の式(2)によって与えられる。このユークリッド距離は、基準座標系の回転、座標系の原点の平行移動、座標軸の鏡映及び座標の順序付けに関し、不変である。そこで、一般性を失わないように、V (506) がクラスタ 1 (505) に属するものとする。次に、クラスタ 1 (505) の共分散行列の固有ベクトルによって定義される基準座標系を考慮し、この基準座標系の原点が重心 1 であるものとする。そうすると、T (501) とクラスタ 1 (505) 内にある V (506) との間の距離

$$\begin{aligned}
 d'^2 &= \sum_{i=1}^k (T_i - V_i)^2 + \sum_{i=k+1}^n (T_i)^2 \\
 &= d_1^2 + d_2^2 \quad (7)
 \end{aligned}$$

【0057】項 d_1 は、部分空間 1 への T (501) の射影 540、すなわち射影 1 と、部分空間 1 への V (506) の射影 V' (507) との間のユークリッド距離である。項 d_2 は、T (501) と、部分空間 1 へのその射影、すなわち射影 1 との間の距離である。換言すれば、項 d_2 は、T (501) と部分空間 1 との間の距離である。かくて、T (501) とベクトル V (506) との間の距離を計算する際に、式(6)に式(7)を代入すると、導かれる近似を束縛することができる。初等幾何学が教えるところによれば、検討中の 3 点、すなわち T (501)、V (506) 及び V' (507) は、一意的な 2 次元部分空間 (平面) を識別する。説明を簡潔にするため、この平面が図5の平面 520 に対応するものとする。そうすると、式(6)で定義されている距離 d は、T (501) 及び V (506) を結合する線分の長さに等しくなり、式(7)で定義されている距離 d' は、T (501) 及び V' (507) を結合する線分の長さに等しくなる。初等幾何学の周知の定理によれば、3 角形の 1 辺の長さは、他の 2 辺の長さの差の絶対値よりも長く、それらの和よりも短い。このことが暗示するのは、式(6)で定義されている距離 d を式(7)で定義されている距離 d' で置換する際に生ずる誤差が、V (506) 及び V' (507) を結合する線分の長さよりも短いか又はこれに等しい、ということであ

は、次のように表すことができる。

【0054】

10 【数6】

$$d^2 = \sum_{i=1}^n (T_i - V_i)^2 \quad (6)$$

【0055】但し、座標 T_i 及び V_i は、基準座標系に対し相対的である。次に、最後の $n - k + 1$ 個の座標をゼロに設定することにより、クラスタ 1 (505) を部分空間 1 に射影する。かくて、T (501) と V (506) を部分空間 1 に射影した点 V' (507) との間の距離は、次のように表すことができる。

20 【0056】

【数7】

る。従って、自乗誤差は、次のように束縛されることになる。

【0058】

30 【数8】

$$error^2 \leq \sum_{i=k+1}^n (V_i)^2 \quad (8)$$

【0059】図6には、次元縮小済みクラスタの階層及びこの階層の最下部にあるクラスタ用の低次元索引を生成するための流れ図が示されている。ステップ 601 で、クラスタ化プロセスは、原データ 602 を入力として受け取り、このデータを複数のデータ・クラスタ 603 に区分するとともに、この区分の詳細に関するクラスタ化情報 604 を生成する。原データ 602 内の各エントリは、式(1)で定義されているようなベクトル属性を含んでいる。クラスタ化アルゴリズムは、当分野で公知のクラスタ化又はベクトル量子化アルゴリズムのうち任意のものを選択することができる。これらのアルゴリズムについては、例えば Leonard Kauffman et al, "Finding Groups in Data", John Wiley & Sons, 1990、又は Yoseph Linde et al, "An Algorithm for Vector Quantizer Design", IEEE Transactions on Communications, Vol. COM-28, No. 1, January 1980, pp. 84-95 を参照されたい。クラスタ化アルゴリズムの学習段階によって生ぜられるクラスタ化情報 604 は、このアルゴ

リズムとともに変化する。すなわち、クラスタ化情報 604 は、クラスタ化アルゴリズムの分類段階によって、新しい任意の、以前に見られなかったサンプルを生成するとともに、各クラスタごとに 1 つの表現ベクトルを生成することを可能にする。クラスタ化情報 604 は、各々が一意的なラベルに関連する、複数のクラスタの重心を含んでいることが好ましい（例えば、前掲の Yoseph Linde et al, "An Algorithm for Vector Quantizer Design" を参照）。ステップ 605 で、順序づけ論理が開始して、後続の動作が各クラスタに対し個別的に且つ順次的に適用されるように、動作の流れを制御する。なお、複数の計算回路が設けられている場合には、順序付け論理 605 をディスパッチング論理によって置き換えることにより、各々が異なるデータ・クラスタについて動作する、複数の計算回路上で同時的な計算を行わせることができる。ステップ 606 で、次元縮小論理 606 は、1 つのデータ・クラスタ 603 を受け取り、次元縮小情報 607 及び次元縮小済みデータ・クラスタ 608 を生成する。ステップ 609 で、終了条件のテストを遂行する（後述）。もし、終了条件が満足されなければ、ステップ 611 で、原データ 602 を次元縮小済みデータ・クラスタ 608 で置換した後、このプロセスをステップ 601 に戻すことにより、ステップ 601 ~ 609 を再帰的に適用することができる。他方、終了条件が満足されると、ステップ 610 で、現クラスタ用の探索可能索引 612 を生成する。もし、ステップ 613 で、既に分析済みのクラスタの数が、ステップ 601 でクラスタ化アルゴリズムによって生成したクラスタ 603 の数と等しいことが分かれば、このプロセスが終了する。さもなければ、このプロセスはステップ 605 に戻る。一般に、クラスタ 603 の数は、ユーザによって選択されるが、当分野では自動式手順も知られている（例えば、Brian Everitt, "Cluster Analysis", Halsted Press, 1973, Chapter 4.2 を参照）。

【0060】ステップ 609 における終了条件のテストは、次のように定義されるデータ・ボリューム $V(X)$ の概念を基礎とすることができる。

【0061】

【数 9】

$$V(X) = \sum_{i=1}^m n_i \quad (9)$$

【0062】但し、 X は複数レコードの集合であり、 n_i は i 番目のレコードであり、和は X の全ての要素についてのものである。もし、これらのレコードが次元縮小ステップ 606 の前に同じサイズ S を有し、そして n が 1 つのクラスタ内にあるレコードの数を表すのであれば、このクラスタのボリュームは $S n$ となる。他方、 S' が次元縮小ステップ 606 の後のレコードのサイズを表すものとすれば、次元縮小後のこのクラスタのボリュームは $S' n$ となる。終了条件は、ボリューム $S n$ 及

び $S' n$ を比較してテストすることができ、 $S n = S' n$ となる場合にこのプロセスが終了する。

【0063】他の実施例では、ステップ 609 における終了条件のテストが存在しないので、このプロセスの再帰的適用は遂行されない。

【0064】図 7 には、ステップ 606 の次元縮小論理が例示されている。図示のように、ステップ 701 で、特異値分解論理は、データ・クラスタ 702 を入力として受け取り、1 つの変換行列 703 及びその複数の固有値 704 を生成する。変換行列 703 の複数の列は、この行列の複数の固有ベクトルである。特異値分解用のアルゴリズムは、当分野では周知である（例えば、R. A. Horn et al, "Matrix Analysis", Cambridge University Press (1985) を参照）。当業者には明らかなように、代替実施例では、ステップ 701 の特異値分解論理を、当分野では周知の主成分分析論理によって置き換えることができる。

【0065】ステップ 705 で、分類論理は、固有値 704 を入力として受け取り、大きさの減少順に分類済みの固有値 706 を生成する。かかる分類論理は、当分野では周知の任意の分類論理とすることができる。ステップ 707 で、選択論理は、所定の選択基準に従って、最大の固有値 708 を保持する処の順序付けられた固有値の部分集合 706 を選択する。かかる選択基準の一例は、その和が変換行列 703 のトレースのユーザ指定パーセントよりも大きい、最小グループの固有値を選択するようなものである。当分野では周知のように、このトレースは、対角線上の要素の和である。この例では、変換行列 703 及び選択済み固有値 708 が、次元縮小情報 607 を構成する。代替的に、固有値の選択を、精度と再現度のトレードオフに基づいて遂行することができる（後述）。

【0066】ステップ 709 で、変換論理は、データ・クラスタ 702 及び変換行列 703 を入力として受け取り、変換行列 703 が指定する変換をデータ・クラスタ 702 の諸要素に適用することによって、変換済みデータ・クラスタ 710 を生成する。ステップ 711 で、選択済み固有値 708 及び変換済みデータ・クラスタ 710 を使用して、次元縮小済みデータ・クラスタ 712 を生成する。推奨実施例では、最小数の次元を保存することにより、次元縮小を遂行している。このように最小数の次元を保存するのは、対応する固有値の集合が全体的な分散のうち少なくとも一定のパーセント（例えば、95%）を占めるようにするためである。

【0067】代替的に、固有値の選択を、精度と再現度のトレードオフに基づいて遂行することもできる。精度と再現度を理解するには、本発明の 1 つの方法によって遂行される探索動作が、（図 10 及び図 11 を参照して後述するように）近似を対象とし得ることに留意する必要がある。ここで、 k を、 N 個の要素を有するデータベ

ースにおいて 1 つのテンプレートに対し最も類似する要素の所望の数であるものとする (k 最近傍)。この探索動作は近似を対象とするものであるから、ユーザは、k よりも多い数の結果が戻されることを要求するのが普通である。n を、戻される結果の数であるものとして、n 個の結果のうち c 個の結果だけが正しい。つまり、c 個の結果が、このテンプレートに対する k 最近傍に含まれている、ということである。精度は、戻された結果と正しい結果の比であって、これを次のように定義する。

$$\text{精度} = c / n$$

再現度は、探索動作によって戻された正しい結果の比であって、これを次のように定義する。

$$\text{再現度} = c / k$$

精度及び再現度はテンプレートの選択に応じて変化するから、それらの予期値は、当該システムの性能の比較的良好な測度となる。かくて、精度及び再現度の両者は、複数のテンプレートの分布にわたって、固定値 n 及び k の関数として取られた予期値 (E) を意味する。

$$\text{精度} = E(c) / n$$

$$\text{再現度} = E(c) / k$$

明らかに、戻される結果の数 n が増大するにつれて、精度が減少するのに対し、再現度は増大する。精度及び再現度の傾向は、単調ではないのが普通である。E(c) は n に依存するから、効率—再現度の曲線は、n のパラメータ関数としてプロットされることが多い。推奨実施例では、要求元が、探索の所望の精度及び許容再現度に関する下限を指定する。その場合、次元縮小論理 (図 6 のステップ 606) は、精度及び再現度に基づいて、次元縮小を遂行する。すなわち、複数の固有値を減少順に順序付けた後、次元縮小論理は、最小の固有値に対応する次元を除去し、そして原トレーニング集合又はユーザが提供するトレーニング集合から無作為に選択したテスト・サンプル集合に基づいて、結果的な精度—再現度関数を推定する。次元縮小論理は、この精度—再現度関数から、所望の再現度が得られるような精度の最大値 n_{max} を導出する。引き続いて、次元縮小論理は、次の最小の固有値を除去して同じ手順を反復し、所望の再現度が得られるような対応する精度を計算する。かかる反復手順が終了するのは、計算済みの精度がユーザ指定のスレッシュホールド値よりも小さくなる場合である。その場

合、次元縮小論理は、かかる終了条件が生ずる直前の反復手順で保存されていた次元だけを保存するのである。

【0068】本発明の他の実施例では、要求元が、所望の再現度の値だけを指定するのに対し、次元縮小論理は、所望の再現度を得るために精度を増大させる場合のコストを推定する。このコストは、2 つの成分を有する。一方の成分は、低次元の空間における距離の計算及び最近傍探索が一層効率的であるという理由で、次元の数に応じて減少する。他方の成分は、所望の再現度を保証するために、保存済み次元の数が減少するにつれて検

索済み結果の数が増大しなければならないという理由で、次元の数に応じて増大する。効率的な方法を使用したとしても、比較的大きな数 n の最近傍を検索することは高価につく。というのは、分析しなければならない探索空間の部分が、所望の結果の数に応じて増大するからである。その場合、次元縮小論理は、全数探索を遂行することにより、ユーザが指定した再現度の値のための探索コストを最小限にする処の、保存すべき次元の数を探索する。

- 10 【0069】クラスタ化及び特異値分解ステップは、終了条件に到達するまで (ステップ 609)、複数のベクトルに対し再帰的に適用することができる (ステップ 601~611)。かかる終了条件の 1 つは、本明細書に記述されているように、各クラスタの次元がもはや縮小不能となるような場合である。オプションとして、R ツリーのような通常の空間索引付け技法を各クラスタに適用することもできる。これらの技法は、次元を最小化したクラスタについては一層効率的である。こうして、高次元ベクトルの集合についての索引生成プロセスが完了する。

【0070】図 8~図 15 を参照して後述するように、本発明は、多次元データのコンパクト表示を使用して効率的な探索を遂行することを特徴としている。当業者には明らかなように、本発明の探索方法は、本明細書に記述されている多次元データの特定のコンパクト表示に限られるものではない。

- 【0071】図 8 には、本発明に従って生成した探索可能な多次元索引 (図 1 の 108 又は図 6 の 612) に基づいて、絶対探索を遂行するための論理の流れが例示されている。この例では、索引を生成するに当たり、クラスタ化及び特異値分解ステップを再帰的に適用していない。絶対探索とは、探索照会 (例えば、探索テンプレート) に正確にマッチする、1 つ以上のレコードを検索するためのプロセスである。図示のように、ステップ 802 で、(クラスタ探索論理とも呼ばれる) 多次元索引エンジン 107 が、探索テンプレート 801 のような指定データを含んでいる 1 つの照会を入力として受け取る。ステップ 802 で、図 6 のステップ 601 で生成したクラスタ化情報 604 を使用して、探索テンプレート 801 が属するクラスタを識別する。ステップ 803 で、図 6 のステップ 606 で生成した次元縮小情報 607 を使用して、探索テンプレート 801 をステップ 802 で識別したクラスタに関連する部分空間に射影することにより、射影済みテンプレート 804 を生成する。ステップ 803 で、クラスタ内探索論理は、図 6 のステップ 610 で生成した探索可能な多次元索引 612 を使用して、射影済みテンプレート 804 についての探索を遂行する。空間索引付け構造が利用不能である場合、最も簡単なクラスタ内の探索機構は、線形走査 (又は線形探索) を遂行するというものである。クラスタの次元が比較的

小さな場合（例えば、10より小さな場合）、Rツリーのような空間索引付け構造は、線形走査よりも良好な効率を有するのが普通である。

【0072】図9には、本発明に従って生成した探索可能な多次元索引（図1の108又は図6の612）に基づいて、絶対探索を遂行するための他の論理の流れが例示されている。この例における索引（108又は612）は、クラスタ化及び次元縮小論理を再帰的に適用して生成されたものである。絶対探索とは、探索照会に正確にマッチする1つ以上のレコードを検索するためのプロセスである。図示のように、ステップ902で、（図8のステップ802におけるクラスタ探索論理に類似する）クラスタ探索論理が、探索テンプレート901のような指定データを含んでいる1つの照会を入力として受け取る。ステップ902で、図6のステップ601で生成したクラスタ化情報604を使用して、探索テンプレート901が属するクラスタを識別する。（図8のステップ803に類似する）ステップ903で、図6のステップ606で生成した次元縮小情報607を使用して、探索テンプレート901をステップ902で識別したクラスタに関連する部分空間に射影することにより、射影済みテンプレート904を生成する。ステップ805で、現クラスタが終端であるか否か、すなわち多次元索引構成プロセスの間に、現クラスタに対しそれ以上の再帰的なクラスタ化及び特異値分解ステップが適用されなかったか否かを決定する。もし、現クラスタが終端でなければ、ステップ907で、探索テンプレート901を射影済みテンプレート904によって置き換えた後、このプロセスはステップ902に戻る。他方、現クラスタが終端であれば、ステップ906で、クラスタ内探索論理は、探索可能な多次元索引612を使用して、射影済みテンプレート904についての探索を遂行する。前述のように、空間索引付け構造が利用不能である場合、最も簡単なクラスタ内探索機構は、線形走査（又は線形探索）を遂行するというものである。クラスタの次元が比較的小さな場合（例えば、10より小さな場合）、Rツリーのような空間索引付け構造は、線形走査よりも良好な効率を有するのが普通である。

【0073】また、本発明は、指定データに最も類似する検索済みのk個の要素のうち最も遠い要素よりも指定データに近いような要素を他のクラスタが保持し得るか否かを評価することを特徴としている。公知のように、クラスタ化情報を使用して複数の区分の境界を再構成することが可能であり、またこれらの境界を使用して1つのクラスタがk最近傍を保持し得るか否か決定することができる。当業者には明らかなように、これらのクラスタ境界は、クラスタ自体の構造に対する簡単な近似である。換言すれば、この境界の数学的形式からは、この境界上の所与の位置にクラスタの要素が存在するか否かを断定することはできない。一例として、データベースが

2つの球形データ・クラスタを保持しており、そしてこれらのクラスタが互いに著しく離れているようなケースを検討する。このケースの妥当な境界は、これらのクラスタの重心を結合する線分に垂直で且つこれらの重心から等距離の超平面となろう。しかしながら、これらのクラスタは互いに著しく離れているから、この境界近傍には如何なるデータ点も存在しないことになる。他のケースでは、この境界は、両クラスタの多数の要素に非常に接近することがあり得る。

10 【0074】図14及び図15を参照して後述するように、本発明は、クラスタ化情報に加えて、各クラスタの実際の形状構造に対する近似階層を計算及び記憶するとともに、かかる近似階層を使用することにより、所与のベクトルからの一定の距離よりも近いような要素を保持し得るクラスタを識別することを特徴としている。

【0075】図10には、本発明に従って生成した探索可能な多次元索引（図6の612）に基づいて、k最近傍探索プロセスを遂行するための論理の流れが例示されている。この例では、索引を生成するに当たり、クラスタ化及び特異値分解ステップを再帰的に適用していない。k最近傍探索とは、探索テンプレートの形式を有する指定データに最も類似する、データベース内のk個のエントリを戻すためのプロセスである。ステップ1002で、所望のマッチの数に等しいk（1000）を使用して、k最近傍集合1009を初期化する。その目的は、k最近傍集合1009が高々k個の要素を保持し、しかも次のステップが開始する前に空であるようにすることである。ステップ1003で、クラスタ探索論理は、探索テンプレート1001のような照会を入力として受け取るとともに、図6のステップ601で生成したクラスタ化情報604を使用して、探索テンプレート1001が属するクラスタを識別する。ステップ1004で、次元縮小情報607を使用して、探索テンプレート1001をステップ1003で識別したクラスタに関連する部分空間に射影する。この射影ステップ1004は、射影済みテンプレート1006及び次元縮小情報1005を生成する。後者の次元縮小情報1005は、射影済みテンプレート1006の直交補空間（探索テンプレート1001及び射影済みテンプレート1006のベクトル差によって定義されるもの）及びこの直交補空間のユークリッド距離を含んでいる。ステップ1007で、クラスタ内探索論理は、次元縮小情報1005及び射影済みテンプレート1006に加え、探索可能な多次元索引612を使用して、k最近傍集合1009を更新することができる。本発明に従って適応可能なクラスタ内探索論理は、公知の最近傍探索方法のうち任意のものとしてすることができる（例えば、公知の方法については、“Nearest Neighbor Pattern Classification”, B. V. Desathary (editor), IEEE Computer Society (1991)を参照）。本発明に従ったクラスタ内探索論理（ス

ステップ 1 0 0 7) は、例えば、射影済みテンプレート 1 0 0 6 と縮小済み次元を有するベクトル空間内にあるクラスタの諸メンバとの間の自乗距離を計算するステップと、その結果を探索テンプレート 1 0 0 1 とこのクラスタの部分空間との間の自乗距離に加えるステップと、その最終的な結果を直交補空間の自乗長さの「和」として定義するステップとを含んでいる。この論理の結果は、次元縮小情報 1 0 0 5 の一部として、ステップ 1 0 0 4 で次のように計算される。

$$\delta^2 \text{ (テンプレート, 要素)} = D^2 \text{ (射影済みテンプレート, 要素)} + \Sigma \parallel \text{直交補空間} \parallel^2.$$

【0 0 7 6】もし、ステップ 1 0 0 7 の開始時に、k 最近傍集合 1 0 0 9 が空であれば、クラスタ内探索論理は、現クラスタ内にある要素の数が k より大きいときは、射影済みテンプレート 1 0 0 6 に最も近い現クラスタの k 個の要素で、さもなければ現クラスタの全ての要素で、k 最近傍集合 1 0 0 9 を充填してこれを更新する。k 最近傍集合 1 0 0 9 の各要素は、対応するミスマッチ索引 δ^2 に関連付けられている。

【0 0 7 7】他方、ステップ 1 0 0 7 の開始時に、k 最近傍集合 1 0 0 9 が空でなければ、クラスタ内探索論理は、k 最近傍集合 1 0 0 9 内の要素に現に関連する最大の索引よりも小さなミスマッチ索引 δ^2 を有するような要素が探索されるときに、k 最近傍集合 1 0 0 9 を更新する。k 最近傍集合 1 0 0 9 を更新するためには、そこから最大のミスマッチ索引 δ^2 を有する要素を除去するとともに、これを新しく探索された要素で置き換えることができる。

【0 0 7 8】もし、k 最近傍集合 1 0 0 9 が k 個よりも少ない要素を保持していれば、不在の要素は、無限の距離にある要素と見なされる。ステップ 1 0 0 8 で、最近傍を保持し得る他の候補クラスタの存否を決定する。このステップは、クラスタ化情報 6 0 4 を入力として受け取り、これからクラスタ境界を決定することができる。もし、探索テンプレート 1 0 0 1 が属さないようなクラスタの境界が、k 最近傍集合 1 0 0 9 の最も遠い要素よりも近ければ、このクラスタは候補クラスタである。もし、候補クラスタが全く存在しなければ、このプロセスは終了し、k 最近傍集合 1 0 0 9 の内容が結果として戻される。さもなければ、このプロセスはステップ 1 0 0 4 に戻り、そこで現クラスタがステップ 1 0 0 8 で識別した候補クラスタとなる。

【0 0 7 9】図 1 1 には、本発明に従って生成した探索可能な多次元索引 (図 6 の 6 1 2) に基づいて、k 最近傍探索プロセスを遂行するための論理の流れが例示されている。この例では、索引を生成するに当たり、クラスタ化及び特異値分解ステップを再帰的に適用している。k 最近傍探索とは、探索テンプレートの形式を有する指定データに最も類似する、データベース内の k 個のエントリを戻すためのプロセスである。ステップ 1 1 0 2

で、k 最近傍集合 1 1 1 1 を空に初期化し、そしてこれが高々 k 個の要素を保持し得るように、所望のマッチの数 k (1 1 0 0) を使用して、k 最近傍集合 1 1 1 1 を初期化する。ステップ 1 1 0 3 で、クラスタ探索論理

は、探索テンプレート 1 1 0 1 を入力として受け取るとともに、図 6 のステップ 6 0 1 で生成したクラスタ化情報 6 0 4 を使用して、探索テンプレート 1 1 0 1 を対応するクラスタに関連付ける。ステップ 1 1 0 4 で、図 6 のステップ 6 0 6 で生成した次元縮小情報 6 0 7 を使用し、探索テンプレート 1 1 0 1 をステップ 1 1 0 3 で識別したクラスタに関連付けられた部分空間に射影する。

この射影ステップ 1 1 0 4 は、射影済みテンプレート 1 1 0 6 及び次元縮小情報 1 1 0 5 を生成する。後者の次元縮小情報 1 1 0 5 は、射影済みテンプレート 1 1 0 6 の直交補空間 (探索テンプレート 1 1 0 1 及び射影済みテンプレート 1 1 0 6 のベクトル差によって定義されるもの) 及びこの直交補空間のユークリッド距離を含んでいることが好ましい。ステップ 1 1 0 7 で、現クラスタが終端であるか否か、すなわち索引を構成する間に、現クラスタに対しそれ以上の再帰的なクラスタ化及び特異値分解ステップが適用されなかったか否かを決定する。

もし、現クラスタが終端でなければ、ステップ 1 1 0 8 で、探索テンプレート 1 1 0 1 を射影済みテンプレート 1 1 0 6 によって置き換えた後、このプロセスはステップ 1 1 0 3 に戻る。さもなければ、ステップ 1 1 0 9 で、クラスタ内探索論理は、次元縮小情報 1 1 0 5 及び射影済みテンプレート 1 1 0 6 に加え、探索可能な多次元索引 6 1 2 を使用して、k 最近傍集合 1 1 1 1 を更新する。本発明に従って適応可能なクラスタ内探索論理

は、当分野で公知の最近傍探索方法のうち任意のものとすることができる (例えば、前掲の "Nearest Neighbor Pattern Classification", B. V. Desathary (editor), IEEE Computer Society (1991) を参照)。本発明に従ったクラスタ内探索論理 (ステップ 1 1 0 9) は、例えば、射影済みテンプレート 1 1 0 6 と縮小済み次元を有するベクトル空間内にある現クラスタの諸メンバとの間の自乗距離を計算するステップと、その結果を探索テンプレート 1 1 0 1 と現クラスタの部分空間との間の自乗距離に加えるステップと、その最終的な結果を直交補空間の自乗長さの「和」として定義するステップとを含んでいる。この論理の結果は、次元縮小情報 1 1 0 5 の一部として、ステップ 1 1 0 4 で次のように計算される。

$$\delta^2 \text{ (テンプレート, 要素)} = D^2 \text{ (射影済みテンプレート, 要素)} + \Sigma \parallel \text{直交補空間} \parallel^2.$$

【0 0 8 0】もし、ステップ 1 1 0 9 の開始時に、k 最近傍集合 1 1 1 1 が空であれば、クラスタ内探索論理は、現クラスタ内にある要素の数が k 個よりも大きいときは、射影済みテンプレート 1 1 0 6 に最も近い現クラスタの k 個の要素で、或いは現クラスタ内にある要素の

数が k に等しいか又はこれより小さいときは、現クラスタの全ての要素で、 k 最近傍集合 1 1 1 1 を充填してこれを更新する。 k 最近傍集合 1 1 1 1 の各要素は、その対応するミスマッチ索引 δ^2 に関連付けられていることが好ましい。

【0081】もし、ステップ 1 1 0 9 の開始時に、 k 最近傍集合 1 1 1 1 が空でなければ、クラスタ内探索論理は、 k 最近傍集合 1 1 1 1 内の要素に現に関連する最大の索引よりも小さなミスマッチ索引 δ^2 を有するような要素が探索されるときに、 k 最近傍集合 1 1 1 1 を更新する。 k 最近傍集合 1 1 1 1 を更新するためには、そこから最大のミスマッチ索引 δ^2 を有する要素を除去するとともに、これを新しく探索された要素で置き換えればよい。

【0082】もし、 k 最近傍集合 1 1 1 1 が k 個よりも少ない要素を保持していれば、不在の要素は、無限の距離にある要素と見なされる。ステップ 1 1 1 0 で、(最初のクラスタ化ステップを適用する前に) 階層の現レベルが、その最上レベルであるか否かを決定する。もし、現レベルが最上レベルであれば、このプロセスは終了し、 k 最近傍集合 1 1 1 1 の内容が結果として戻される。他方、現レベルが最上レベルでなければ、ステップ 1 1 1 4 で、現レベルにおける候補クラスタ、すなわち k 最近傍のうち幾つかのものを保持し得るクラスタを探索する。この探索は、次元縮小情報 1 1 0 5 及びクラスタ化情報 6 0 4 を使用して遂行される。ステップ 1 1 1 4 で、クラスタ化情報 6 0 4 を使用して、クラスタ境界を決定する。もし、探索テンプレート 1 1 0 1 が属さないようなクラスタの境界が、 k 最近傍集合 1 1 1 1 の最も遠い要素よりも近ければ、このクラスタは候補クラスタである。もし、候補クラスタが全く存在しなければ、ステップ 1 1 1 3 で、現レベルを階層の先行レベルに設定し、次元縮小情報 1 1 0 5 を更新した後、このプロセスはステップ 1 1 1 0 に戻る。他方、候補クラスタが存在すれば、ステップ 1 1 1 5 で、テンプレートを候補クラスタに射影して射影済みテンプレート 1 1 0 6 を更新するとともに、次元縮小情報 1 1 0 5 を更新する。次いで、このプロセスはステップ 1 1 0 7 に戻る。

【0083】図 1 2 の (a) ~ (c) は、類似性のみに基づくクラスタ化技法の結果を、データの局所構造に適応させたアルゴリズムを使用したクラスタ化の結果と比較している。類似性のみに基づくクラスタ化技法は、例えば、各クラスタの諸要素と対応する重心との間のユークリッド距離を最小化することを基礎としている (前掲の Linde et al, "An Algorithm for Vector Quantizer Design" を参照)。詳述すると、図 1 2 の (a) は、基準座標系 1 2 0 1 及び複数ベクトルの集合 1 2 0 2 を示している。もし、各クラスタの諸要素と対応する重心との間のユークリッド距離を最小化することを基礎とするクラスタ化技法が使用されているのであれば、その可

能な結果は、図 1 2 の (b) に示されている通りである。すなわち、ベクトル集合 1 2 0 2 は、超平面 1 2 0 3 及び 1 2 0 4 によって、クラスタ 1 (1 2 0 5)、クラスタ 2 (1 2 0 6) 及びクラスタ 3 (1 2 0 7) から成る、3 つのクラスタに区分される。かかる結果的なクラスタは、互いに類似するベクトルを保持しているが、データの局所構造を捕捉しないから、準一最適の次元縮小を与える。図 1 2 の (c) は、データの局所構造に適応させたアルゴリズムを使用したクラスタ化の結果を示している。このクラスタ化から得られる 3 つのクラスタ、すなわちクラスタ 1 (1 2 0 8)、クラスタ 2 (1 2 0 9) 及びクラスタ 3 (1 2 1 0) は、データの局所構造を良好に捕捉するから、独立の次元縮小に対し適応させることができる。

【0084】図 1 3 には、データの局所構造に適応させたクラスタ化アルゴリズムが例示されている。ステップ 1 3 0 2 で、クラスタ化すべきベクトル集合 1 3 0 1 及びクラスタの所望の数 NC を使用して、重心の初期値 1 3 0 3 を選択する。推奨実施例では、置換を伴わない公知のサンプリング技法を使用して、クラスタの所望の数 NC の各々ごとに、ベクトル集合 1 3 0 1 の 1 つの要素を無作為に選択する。ステップ 1 3 0 4 で、例えばユークリッド距離に基づく任意の公知方法を使用して、第 1 のクラスタ集合を生成する。その結果、複数のサンプルが NC 個のクラスタに分割される。ステップ 1 3 0 6 で、 NC 個のクラスタの各々の重心を、例えばこのクラスタ内にあるベクトルの平均として、計算する。ステップ 1 3 0 8 で、特異値分解論理 (図 7 のステップ 7 0 1) を使用して、クラスタ 1 3 0 5 の固有値及び固有ベクトル 1 3 0 9 を計算することができる。ステップ 1 3 1 0 で、重心情報 1 3 0 7、固有ベクトル及び固有値 1 3 0 9 を使用して、各クラスタごとに異なる距離関数を生成する。特定のクラスタの距離関数は、例えば、固有ベクトルによって定義された回転済み空間内のユークリッド距離を、この固有ベクトルの平方根に等しい重みで加重したものである。

【0085】ステップ 1 3 1 2 ~ 1 3 1 4 から成るループは、新しいクラスタを生成する。ステップ 1 3 1 2 で、制御論理は、ベクトル集合 1 3 0 1 内の全てのベクトルにわたって、ステップ 1 3 1 3 及び 1 3 1 4 を反復する。ステップ 1 3 1 3 で、距離関数 1 3 1 1 を使用して、選択したベクトルとクラスタの重心の各々との間の距離を計算する。ステップ 1 3 1 4 で、このベクトルを最も近いクラスタに割り当てることにより、クラスタ 1 3 0 5 を更新する。ステップ 1 3 1 5 で、終了条件に到達していれば、このプロセスを終了する。さもなければ、ステップ 1 3 0 6 で、このプロセスを継続する。推奨実施例では、後続する 2 つの反復の間にクラスタの構成が変化していないことを条件として、このプロセスを終了する。

【0086】図14には、3次元空間における複雑な表面1401と、3次元の4進ツリーに基づく2つの逐次近似1402及び1403が例示されている。かかる逐次近似は、例えば H. Samet, "Region Representation Quadtree from Boundary Codes", Comm. ACM23-3, pp.163-170 (March 1980) の教示に従っている。第1の近似1402は、極小外接ボックスである。これに対し、第2の近似1403は、4進ツリー生成の第2のステップであり、そこでは各次元の midpoint において外接ボックスを分割するとともに、表面に交差する超方形のみを保存することによって、この外接ボックスが8個の超方形に分割されている。

【0087】推奨実施例では、近似の階層を、k 次元の4進ツリーとして生成する。かかる近似の階層を生成するための本発明の方法は、例えば、クラスタの形状に対する0次近似に対応するようなクラスタ境界を生成するステップと、極小外接ボックスにより前記各クラスタの凸包を近似して、前記各クラスタの形状に対する1次近似を生成するステップと、各次元の midpoint において前記外接ボックスを切断することにより、前記外接ボックスを 2^k 個の超方形に区分するステップと、点を保持するような前記超方形のみを保存して、前記各クラスタの形状に対する2次近似を生成するステップと、前記保存した超方形の各々ごとに最後の2つのステップを反復して、前記各クラスタの形状に対する3次、4次、 \dots 、n 次近似を逐次に生成するステップとを含んでいる。

【0088】図15には、クラスタの形状に対する逐次近似を使用して、所与のデータ点から一定距離よりも近い要素を保持し得るクラスタを識別するための論理の流れが例示されている。1つの実施例では、クラスタの形状は、その凸包である。他の実施例では、クラスタの形状は、全ての点を囲んでいるような連結済みの弾性表面である。この論理は、例えば図10のステップ1008において、候補クラスタを探索するのに使用することができる。図15を参照して説明すると、ステップ1502で、形状近似の階層を有するクラスタの原集合1501をこのプロセスに入力するとともに、候補集合1505を原集合1501に初期化する。ステップ1506で、現近似を0次-形状近似に設定して、他の初期化ステップを遂行する。推奨実施例では、諸クラスタの0次-形状近似は、かかるクラスタを生成するために使用したクラスタ化アルゴリズムの決定領域によって与えられる。ステップ1507で、クラスタの形状の現近似とデータ点1503との間の距離を計算する。候補クラスタよりも遠い全てのクラスタを廃棄して、クラスタの保存集合を生成する。ステップ1509で、階層内に良好な近似が存在するか否かを決定する。もし、良好な近似が存在しなければ、結果集合1512を現保存集合150

8に等しく設定して、このプロセスを終了する。さもなければ、ステップ1510で、候補集合1505を現保存集合1508に等しく設定し、現形状近似を階層内の良好な近似に設定した後、このプロセスはステップ1507に戻る。

【図面の簡単な説明】

【図1】本発明に従ったネットワーク化クライアント／サーバ・システムを例示するブロック図である。

【図2】複数データ点の分布及びクラスタ化を遂行した後の次元縮小を直感的に例示する図である。

【図3】3次元空間内の3点を2次元部分空間に射影するに際し、これらの3点のうち任意の2点間の相対距離を維持するようにした一例を示す図である。

【図4】3次元空間内の3点を2次元部分空間に射影するに際し、相対距離のランクに影響を及ぼすようにした一例を示す図である。

【図5】元の空間及び射影済みの部分空間にある点間の距離の計算の一例を示す図である。

【図6】データベース内のデータから多次元索引を生成するための論理を例示する図である。

【図7】データの次元縮小を遂行するための論理の流れを例示する図である。

【図8】クラスタ化及び特異値分解ステップを再帰的に適用しないで生成した索引を使用して、絶対探索を遂行するための論理の流れを例示する図である。

【図9】クラスタ化及び特異値分解ステップを再帰的に適用して生成した索引を使用して、絶対探索を遂行するための論理の流れを例示する図である。

【図10】クラスタ化及び特異値分解ステップを再帰的に適用しないで生成した索引を使用して、k 最近傍探索を遂行するための論理の流れを例示する図である。

【図11】クラスタ化及び特異値分解ステップを再帰的に適用して生成した索引を使用して、k 最近傍探索を遂行するための論理の流れを例示する図である。

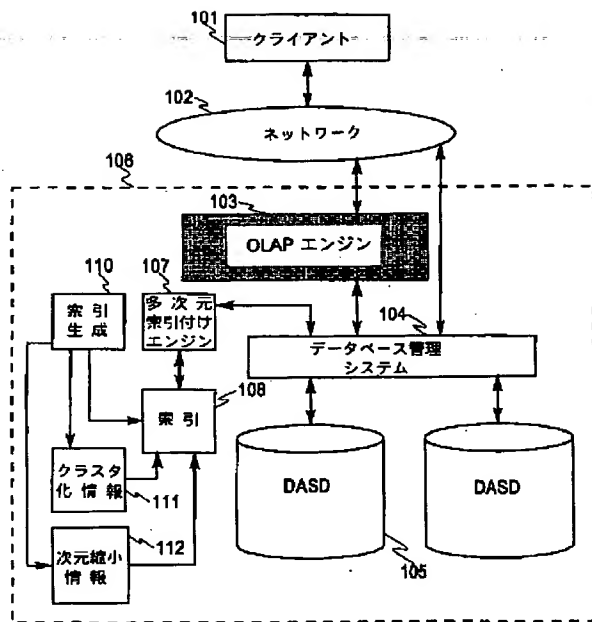
【図12】3次元空間内にあるデータ、並びにユークリッド距離に基づくクラスタ化技法の結果及びデータの局所構造に適応させたクラスタ化技法の結果の比較を例示する図である。

【図13】データの局所構造に適応させたクラスタ化技法の論理の流れを例示する図である。

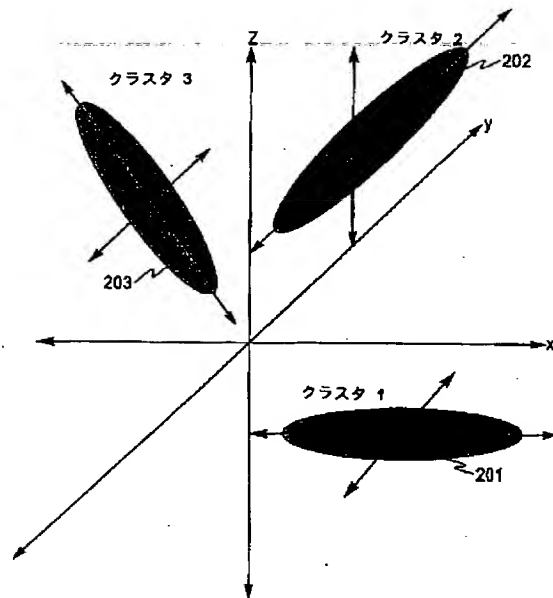
【図14】3次元空間における複雑な超平面と、3次元4進ツリーの生成アルゴリズムを使用して生成した2つの逐次近似を例示する図である。

【図15】クラスタの形状の逐次近似を使用して、所与のベクトルから一定距離よりも近い要素を保持し得るクラスタを決定するための論理の流れを例示する図である。

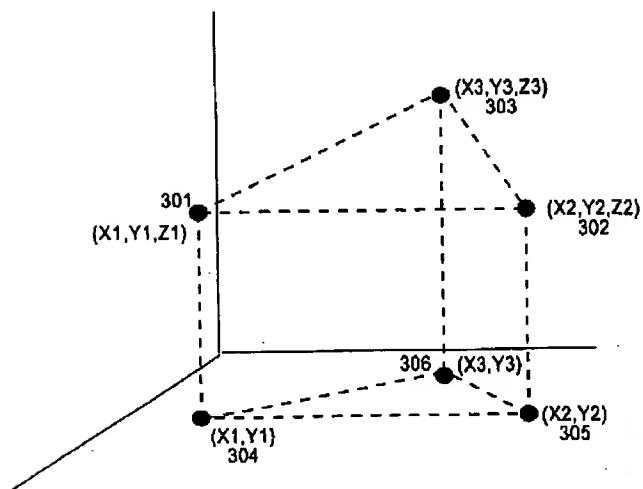
【図 1】



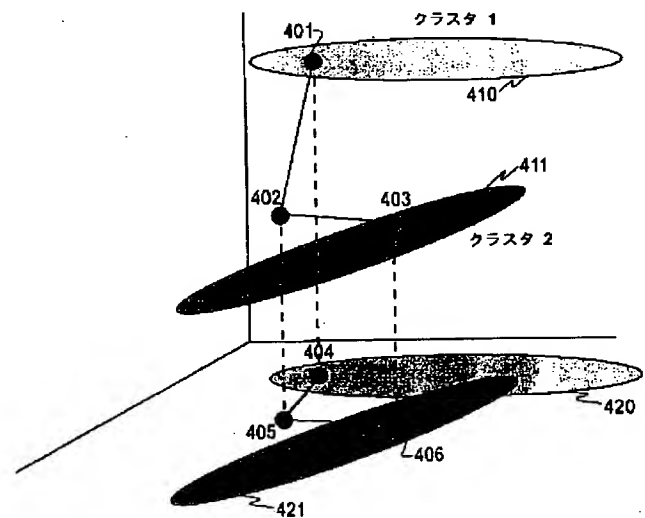
【図 2】



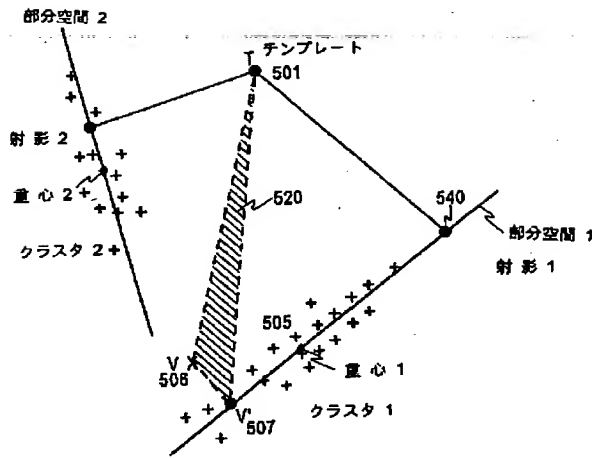
【図 3】



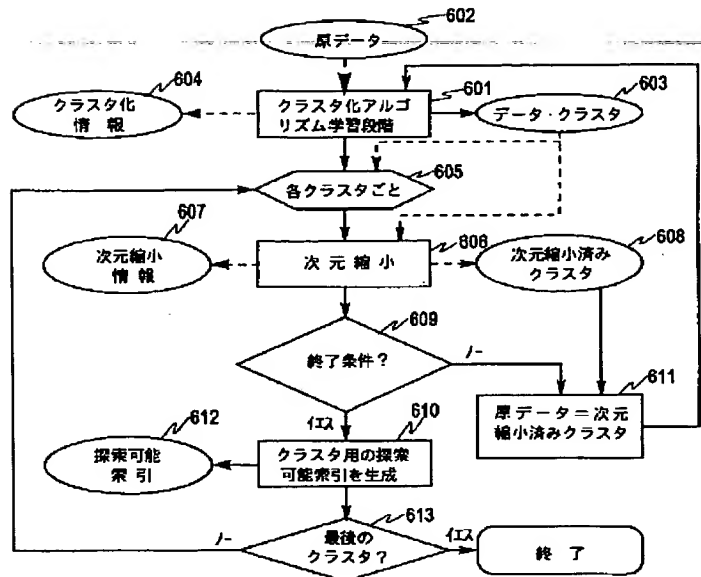
【図 4】



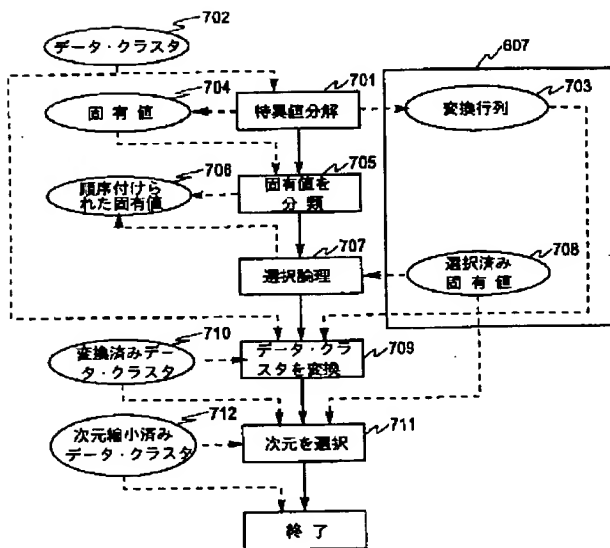
【図 5】



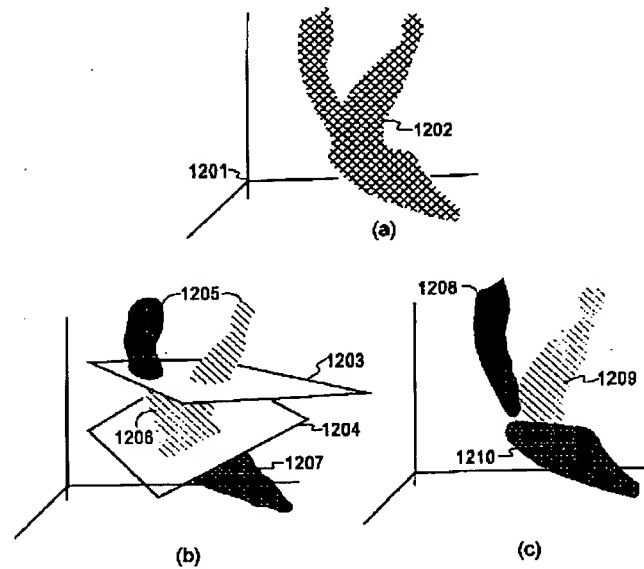
【図 6】



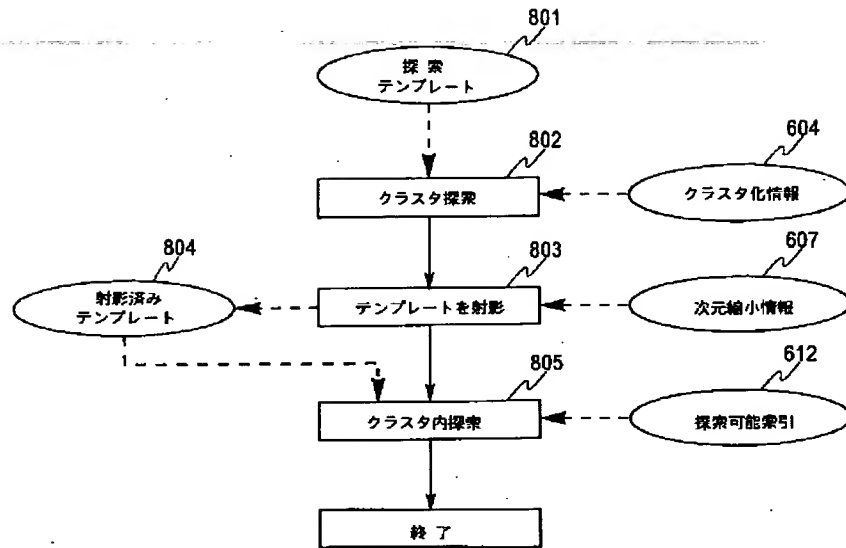
【図 7】



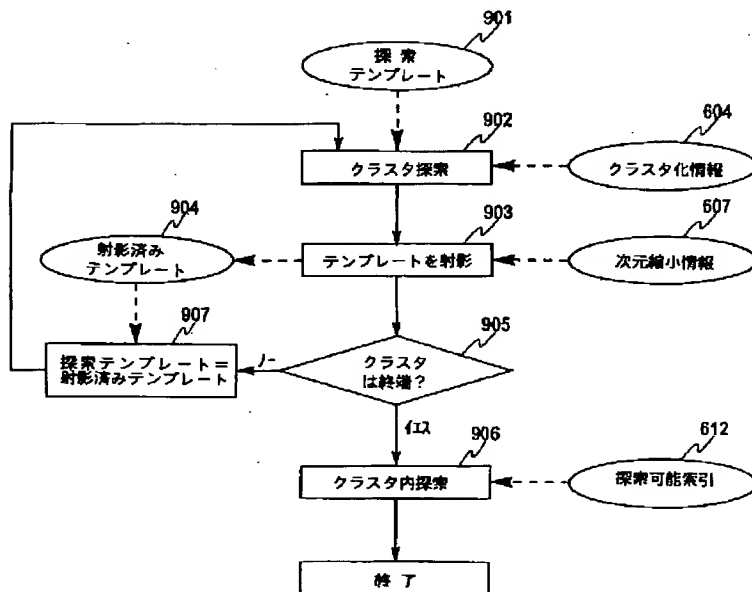
【図 12】



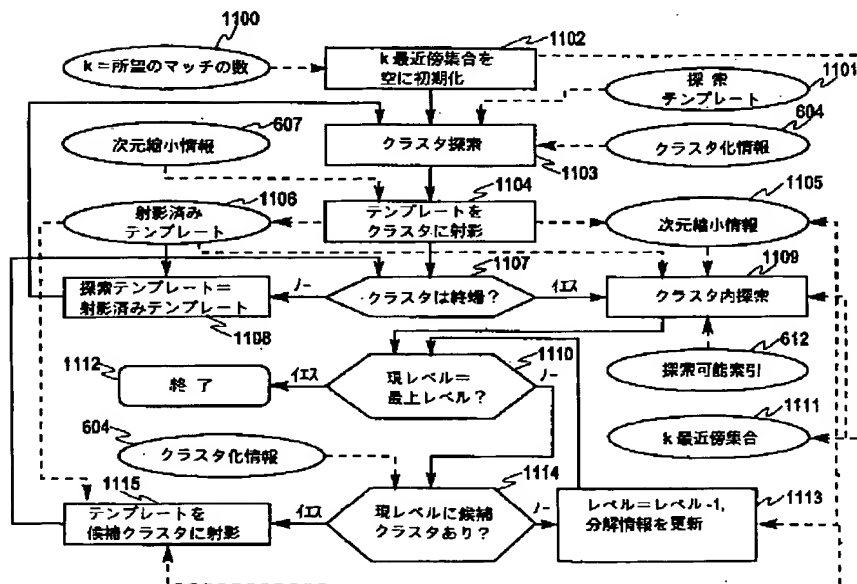
【図 8】



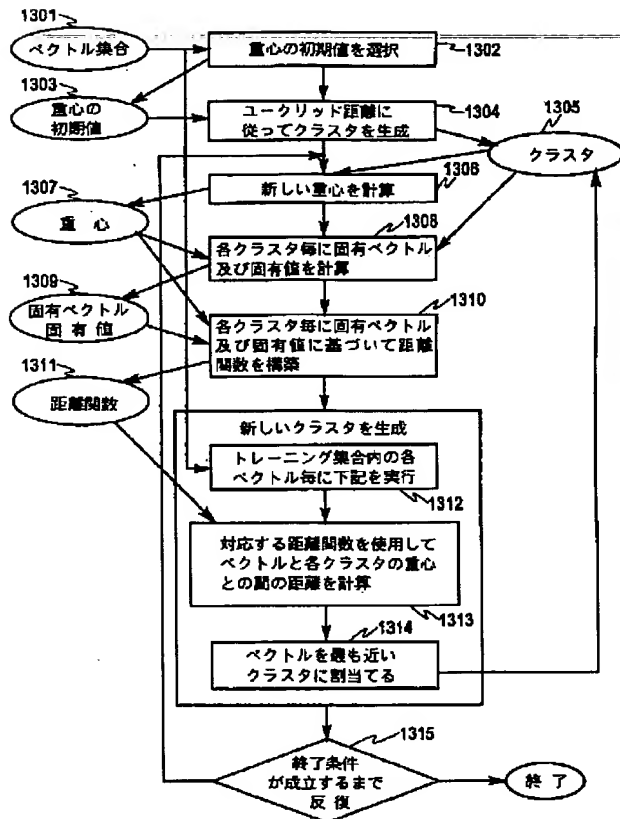
【図 9】



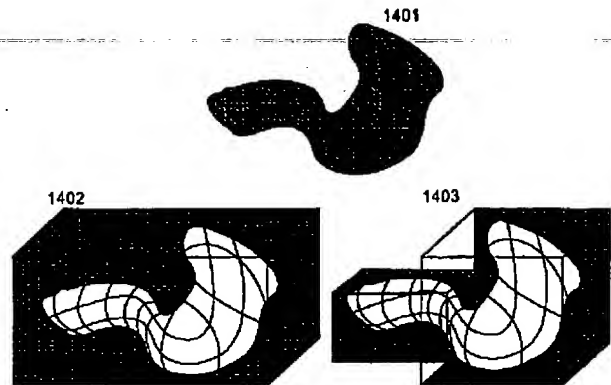
【図 1 1】



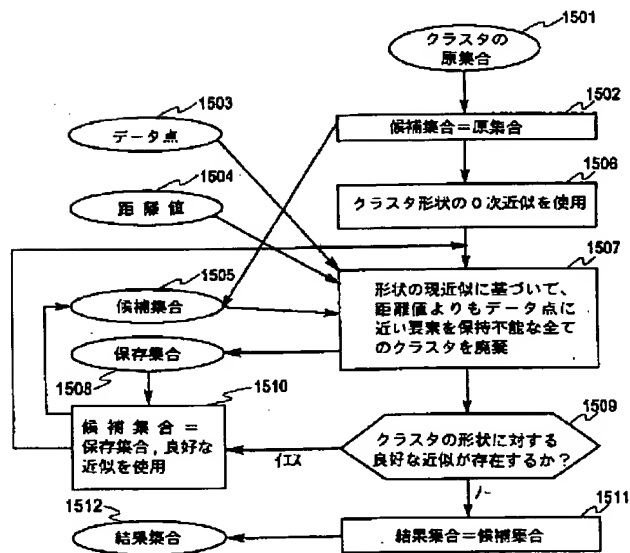
【図 1 3】



【図 1 4】



【図 1 5】



フロントページの続き

(72)発明者 チュン－シェン・リ
アメリカ合衆国10562、ニューヨーク州オ
ッシニング、クロトン・アヴェニュー 50
、アパートメント・2・シー

(72)発明者 アレキサンダー・トマジアン
アメリカ合衆国10570、ニューヨーク州ブ
レザントヴィル、メドウブルック・ロード
17